

April 1991

LIDS-R-2030

LARGE DEVIATIONS AND APPLICATIONS:  
THE FINITE DIMENSIONAL CASE

A. Dembo  
Department of Mathematics and Department of Statistics  
Stanford University  
Stanford, CA 94305, USA

O. Zeitouni  
Department of Electrical Engineering, Technion  
Haifa 32000, Israel

---

This research was supported by the Air Force Office of Scientific Research under contract AFOSR-89-0276B and by the Army Research Office under contract DAAL03-86-K-0171.

Report Documentation Page				Form Approved OMB No. 0704-0188	
Public reporting burden for the collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to a penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.					
1. REPORT DATE <b>APR 1991</b>		2. REPORT TYPE		3. DATES COVERED <b>00-04-1991 to 00-04-1991</b>	
4. TITLE AND SUBTITLE <b>Large Deviations and Applications: The Finite Dimensional Case</b>				5a. CONTRACT NUMBER	
				5b. GRANT NUMBER	
				5c. PROGRAM ELEMENT NUMBER	
6. AUTHOR(S)				5d. PROJECT NUMBER	
				5e. TASK NUMBER	
				5f. WORK UNIT NUMBER	
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) <b>Massachusetts Institute of Technology, Laboratory for Information and Decision Systems, 77 Massachusetts Avenue, Cambridge, MA, 02139-4307</b>				8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)				10. SPONSOR/MONITOR'S ACRONYM(S)	
				11. SPONSOR/MONITOR'S REPORT NUMBER(S)	
12. DISTRIBUTION/AVAILABILITY STATEMENT <b>Approved for public release; distribution unlimited</b>					
13. SUPPLEMENTARY NOTES					
14. ABSTRACT					
15. SUBJECT TERMS					
16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT	18. NUMBER OF PAGES <b>99</b>	19a. NAME OF RESPONSIBLE PERSON
a. REPORT <b>unclassified</b>	b. ABSTRACT <b>unclassified</b>	c. THIS PAGE <b>unclassified</b>			

## Preface and Warning

These notes, which form the first chapter of a forthcoming book, are intended to serve as lecture notes on the topic of large deviations and applications for students whose background and interests are in applications which involve finite dimensional spaces. Although narrow in their scope, these notes present a good deal of the methods available for more general situations. A glaring omission is the method of sub-additivity, which will be discussed in another chapter in the book. Another deficiency of these notes is the sketchy bibliography and historical. We hope to correct this in the book.

The funny looking ?? appearing in various places mean references to later chapters. Please disregard those.

**Acknowledgments** Special thanks go to Dan Stroock who taught one of us (O.Z.) large deviations theory and provided many comments and useful insights. This course was taught at Stanford and the Technion. The comments of people who sat in these courses, and in particular the comments of Yuval Peres, contributed much to correct mistakes and omissions. We wish to thank Sam Karlin for suggesting the application of rare long segments in random walks (section 2.5). One of us (A.D.) thanks Tom Cover and Joy Thomas for a preprint of their forthcoming book which influenced our treatment of sections 2.1.1 and 2.7.

# Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
1.1	The Large Deviations Principle . . . . .	4
1.2	Major large deviations techniques . . . . .	8
<b>2</b>	<b>Large Deviations principles for finite dimensional spaces</b>	<b>10</b>
2.1	Combinatorial Techniques for finite alphabet . . . . .	10
2.1.1	The method of types and Sanov's theorem . . . . .	11
2.1.2	Cramer's theorem for finite alphabets in $\mathbb{R}^1$ . . . . .	17
2.1.3	Large deviations for sampling without replacements . . . . .	19
2.2	Cramer's Theorem in $\mathbb{R}^1$ . . . . .	24
2.3	A general large deviations principle in $\mathbb{R}^d$ . . . . .	34
2.4	Large deviations of Markov chains over finite alphabets . . . . .	50
2.4.1	Cramer's theorem for Markov additive processes . . . . .	51
2.4.2	Sanov's theorem for the empirical measure of Markov chains . . . . .	53
2.4.3	Sanov's theorem for the pair empirical measure of Markov chains . . . . .	56
2.5	Long rare segments in random walks . . . . .	60
2.6	The Gibbs conditioning principle in finite alphabet . . . . .	64
2.7	The hypothesis test problem . . . . .	68
2.8	Generalized maximum likelihood for finite alphabets . . . . .	74

2.9	Rate distortion theory for stationary and ergodic sources . . . . .	79
2.10	Refinements of large deviations statements in $\mathbb{R}^d$ . . . . .	88
<b>3</b>	<b>Historical notes and references</b>	<b>94</b>
3.1	Chapter 1 . . . . .	94
	<b>Bibliography</b>	<b>98</b>

# Chapter 1

## Introduction

### 1.1 The Large Deviations Principle

A *large deviations principle* characterizes the limiting behavior (as  $\epsilon \rightarrow 0$ ) of a family of probability measures  $\mu_\epsilon$  on  $(\mathcal{X}, \mathcal{B}_\mathcal{X})$  in terms of a *rate function*. This characterization is via ‘tight’ asymptotic upper and lower exponential bounds on the values that  $\mu_\epsilon$  assigns to close and open subsets of  $\mathcal{X}$  (respectively). For that purpose  $\mathcal{X}$  should be a topological space so that open and closed subsets of  $\mathcal{X}$  are well defined concepts and as usual only measurable sets (i.e., elements of  $\mathcal{B}_\mathcal{X}$ , the Borel sigma field on  $\mathcal{X}$ ) are of interest.

**Definitions:** A **rate function**  $I$  is any mapping  $I : \mathcal{X} \rightarrow [0, \infty]$  such that for any  $\alpha \in [0, \infty)$  the level set  $\Psi_I(\alpha) \triangleq \{x : I(x) \leq \alpha\}$  is a closed subset of  $\mathcal{X}$ . A **good rate function** is a rate function for which all the level sets  $\Psi_I(\alpha)$  are compact subsets of  $\mathcal{X}$ .

Alternatively, a rate function is any non-negative, lower semicontinuous function on  $\mathcal{X}$  (for the definitions and a proof of this fact see Appendix ??). Throughout, let  $\mathcal{D}_I$  denote the set of points in  $\mathcal{X}$  of *finite rate*, namely  $\mathcal{D}_I \triangleq \{x : I(x) < \infty\}$ .

**Definition:** We say that  $\mu_\epsilon$  satisfies the **large deviations principle with rate function**  $I$  if, for all  $\Gamma \in \mathcal{B}_\mathcal{X}$ ,

$$-\inf_{x \in \Gamma^\circ} I(x) \leq \liminf_{\epsilon \rightarrow 0} \epsilon \log \mu_\epsilon(\Gamma) \leq \limsup_{\epsilon \rightarrow 0} \epsilon \log \mu_\epsilon(\Gamma) \leq -\inf_{x \in \bar{\Gamma}} I(x) \quad (1.1.1)$$

where  $\bar{\Gamma}$  is the closure of  $\Gamma$ ,  $\Gamma^\circ$  is the interior of  $\Gamma$  and the infimum of  $I$  over an empty set is

interpreted throughout as  $\infty$ .

**Corollary 1.1.1** *When  $\mu_\epsilon$  satisfies a large deviations principle with rate function  $I$  and  $\Gamma \in \mathcal{B}_X$  is such that*

$$\inf_{x \in \Gamma^\circ} I(x) = \inf_{x \in \bar{\Gamma}} I(x) \triangleq I_\Gamma \quad (1.1.2)$$

*then*

$$\lim_{\epsilon \rightarrow 0} \epsilon \log \mu_\epsilon(\Gamma) = -I_\Gamma \quad (1.1.3)$$

**Remark:** When  $I$  is a good rate function and  $\bar{\Gamma} \cap \mathcal{D}_I$  is non-empty then there exists at least one point  $x^* \in \bar{\Gamma}$  where  $\inf_{x \in \bar{\Gamma}} I(x)$  is achieved.

A set  $\Gamma$  which satisfies (1.1.2) is called an  *$I$  Continuity Set*. In general, a large deviations principle implies a precise limit as in (1.1.3) only for  $I$  continuity sets.

In particular, points are typically not open subsets of  $\mathcal{X}$  so a large deviations principle *does not* result with an asymptotic exponential estimate on the probability that  $\mu_\epsilon$  assigns to each point in  $\mathcal{X}$ . Better results may well be derived on a case by case basis for specific families of measures  $\mu_\epsilon$  and particular sets. While such results do not fall within our definition of a large deviations principle few illustrative examples are included in this book (see Sections 2.1, 2.6, 2.10).

An alternative formulation of the large deviations principle is as follows:

(a). For any closed set  $F \in \mathcal{B}_X$ ,

$$\limsup_{\epsilon \rightarrow 0} \epsilon \log \mu_\epsilon(F) \leq - \inf_{x \in F} I(x) \quad (1.1.4)$$

(b). For any open set  $G \in \mathcal{B}_X$ ,

$$\liminf_{\epsilon \rightarrow 0} \epsilon \log \mu_\epsilon(G) \geq - \inf_{x \in G} I(x) \quad (1.1.5)$$

The inequality (1.1.4) is also called the large deviations upper bound while (1.1.5) is called the large deviations lower bound. Thus, the large deviations principle corresponds to the scenario in which both bounds hold with the same rate function.

A few observations are now in place. First, observe that since  $\mu_\epsilon(\mathcal{X}) = 1$  for any  $\epsilon$  it is necessary that  $\inf_{x \in \mathcal{X}} I(x) = 0$  for the upper bound to hold (and when  $I$  is a good rate function then there

exists at least one point  $x$  in which  $I(x) = 0$ ). Next, observe that the upper bound trivially holds whenever  $\inf_{x \in F} I(x) = 0$  while the lower bound trivially holds whenever  $\inf_{x \in G} I(x) = \infty$ . This leads to another alternative formulation of the large deviations principle which is useful when trying to prove such a principle.

(a). For any closed set  $F \in \mathcal{B}_{\mathcal{X}}$  contained in the complement of  $\Psi_I(\alpha)$ ,

$$\limsup_{\epsilon \rightarrow 0} \epsilon \log \mu_{\epsilon}(F) \leq -\alpha \quad (1.1.6)$$

(b). For any  $x \in \mathcal{D}_I$  and any open neighborhood  $G \in \mathcal{B}_{\mathcal{X}}$  of  $x$  in  $\mathcal{X}$ ,

$$\liminf_{\epsilon \rightarrow 0} \epsilon \log \mu_{\epsilon}(G) \geq -I(x) \quad (1.1.7)$$

The inequality (1.1.7) reveals the local nature of the lower bound which should only be proved for “small open sets”. Moreover, the following indicates that the upper bound may be proved first for an approximate  $I$  functional:

**Definition:** For any  $\delta > 0$  let  $I^{\delta}(x) = I(x) - \delta$  when  $x \in \mathcal{D}_I$  and  $I^{\delta}(x) \triangleq \frac{1}{\delta}$  when  $x \in \mathcal{D}_I^c$ .

Since for any set  $F$

$$\lim_{\delta \rightarrow 0} \inf_{x \in F} I^{\delta}(x) = \inf_{x \in F} I(x), \quad (1.1.8)$$

it suffices to prove that for any  $\delta > 0$  and for any closed set  $F$

$$\limsup_{\epsilon \rightarrow 0} \epsilon \log \mu_{\epsilon}(F) \leq - \inf_{x \in F} I^{\delta}(x) \quad (1.1.9)$$

in order to conclude that the upper bound (1.1.4) holds.

Actually, a common technique for proving the existence of a large deviations principle is by implicitly defining

$$I(x) \triangleq \sup_{\{A \in \mathcal{C}^{\circ} : x \in A\}} \mathcal{L}_A \quad (1.1.10)$$

where

$$\mathcal{L}_A \triangleq - \liminf_{\epsilon \rightarrow 0} \epsilon \log \mu_{\epsilon}(A), \quad (1.1.11)$$

and  $\mathcal{C}^{\circ} \subset \mathcal{B}_{\mathcal{X}}$  is any basis for the topology of  $\mathcal{X}$  (namely, any open set  $G$  is the union of sets from  $\mathcal{C}^{\circ}$ ).



This definition results with a rate function  $I$  for which the lower bound (1.1.5) holds. Moreover, when  $\lim_{\epsilon \rightarrow 0} \epsilon \log \mu_\epsilon(A)$  exists for any  $A \in \mathcal{C}^0$ , then the upper bound (1.1.4) also holds for any compact  $F \in \mathcal{B}_X$  (these statements are proved in Section ??).

This approach in which the upper bound is first proved directly for compact sets and then its validity extended to all closed sets serves as a motivation for the following definition.

**Definition:** A sequence of probability measures  $\mu_\epsilon$  satisfies **weakly** the large deviations principle with rate function  $I$  if the upper bound (1.1.4) holds for  $F$  **compact** and the lower bound (1.1.5) holds.

The following auxiliary *exponential tightness* property suffices for extending the *weak* large deviations principle to a full large deviations principle with a good rate function (this is shown in the sequel).

**Definition:** A family of measures  $\mu_\epsilon$  is **exponentially tight** if for any  $\alpha < \infty$  there exists a compact set  $K_\alpha \in \mathcal{B}_X$  such that

$$\limsup_{\epsilon \rightarrow 0} \epsilon \log \mu_\epsilon(K_\alpha^c) < -\alpha \quad (1.1.12)$$

where  $K_\alpha^c$  denotes the complement of the set  $K_\alpha$ .

**Remark:** The measures  $\mu_\epsilon$  may satisfy a large deviations principle with a good rate function without being exponentially tight. Beware of this common logical mistake.

The exponential tightness (and the alternative statement of the upper bound (1.1.9)) are applied in the following lemma for strengthening a weak large deviations principle.

**Lemma 1.1.1** *Let  $\mu_\epsilon$  be an exponentially tight family of probability measures.*

- (a). *If the upper bound (1.1.4) holds for all compact sets then it also holds for all closed sets.*
- (b). *If the lower bound (1.1.5) holds for all open sets then  $I(\cdot)$  is a good rate function.*

Thus, when an exponentially tight family of measures satisfies weakly the large deviations principle with a rate function  $I(\cdot)$  then  $I$  is a good rate function and a full large deviations principle holds.

**Proof:** (a). Consider an arbitrary closed set  $F$ . All we need is to establish (1.1.6) whenever  $F \subset \Psi_I(\alpha)^c$ . Fix any such  $\alpha < \infty$ . Clearly,

$$\mu_\epsilon(F) \leq \mu_\epsilon(F \cap K_\alpha) + \mu_\epsilon(K_\alpha^c),$$

where  $K_\alpha$  is the compact set in (1.1.12). Part (a) of the lemma now follows by the inequality (1.1.12) and the upper bound (1.1.4) for the compact set  $F \cap K_\alpha$  (note that  $F \cap K_\alpha \subset \Psi_I(\alpha)^c$ , so  $\inf_{x \in F \cap K_\alpha} I(x) \geq \alpha$ ).

(b). Applying the lower bound (1.1.5) to the open set  $K_\alpha^c$ , one concludes by (1.1.12) that  $\inf_{x \in K_\alpha^c} I(x) > \alpha$ . Therefore,  $\Psi_I(\alpha) \subseteq K_\alpha$ , yielding the compactness of the closed level set  $\Psi_I(\alpha)$ . As this argument holds for any  $\alpha < \infty$ , it follows that  $I(\cdot)$  is a good rate function.  $\square$

A countable family of measures  $\mu_n$  is considered in many cases (for example when  $\mu_n$  is the law governing the empirical mean of  $n$  random variables). Then, a large deviations principle corresponds to the statement

$$-\inf_{x \in \Gamma^o} I(x) \leq \liminf_{n \rightarrow \infty} a_n \log \mu_n(\Gamma) \leq \limsup_{n \rightarrow \infty} a_n \log \mu_n(\Gamma) \leq -\inf_{x \in \bar{\Gamma}} I(x), \quad (1.1.13)$$

for some sequence  $a_n \rightarrow 0$ . Note that here  $a_n$  replaces  $\epsilon$  of (1.1.1) and similarly the statements (1.1.4)-(1.1.7) may be appropriately modified.

## 1.2 Major large deviations techniques

The major disadvantage of the very general approach outlined above is that it does not really reveal what the rate function values are. However, since the rate function associated with a large deviations principle is, under mild conditions, unique (as proved in Section ??) the following two step method is quite useful:

- (a). Prove the existence of a large deviations principle.
- (b). Verify certain properties of the rate function (typically convexity and/or smoothness) and from them deduce a more convenient (explicit) characterization of this function.

Indeed this method prevails to a certain degree throughout Chapters ??-?? of this book.

Nevertheless, when  $\mathcal{X}$  is a subset of a finite dimensional vector space (specifically when  $\mathcal{X} \subset \mathbb{R}^d$  for some  $d < \infty$ ) then one can typically prove directly the existence of a large deviations principle with an **explicit** rate function. Chapter 2 of this book which is dedicated to this class of problems serves at least three purposes:

- (a). To illustrate the sharpest possible results and emphasize both the reason for the existence of a large deviations principle and the types of rate functions one expects to find.
- (b). To demonstrate the different methods for proving explicit large deviations statements in simple scenarios when one's eyesight is not yet obscured by various technical details.
- (c). To present an important class of interesting results and their applications while requiring relatively little mathematical background.

## Chapter 2

# Large Deviations principles for finite dimensional spaces

### 2.1 Combinatorial Techniques for finite alphabet

Throughout this section all random variables assume values in the finite set  $\Sigma = \{a_1, a_2, \dots, a_{|\Sigma|}\}$  ( $\Sigma$  is also sometimes called the underlying alphabet). Combinatorial methods are then applicable for deriving large deviations principles for empirical measures (see Sections 2.1.1 and 2.1.3 below) and for empirical means (see Section 2.1.2 below). While the scope of these methods is limited to finite alphabets, they illustrate the results one can hope to (and should indeed) derive for more abstract alphabets. Some of these results are actually direct consequences of the large deviations principles derived below via the combinatorial method. For example, in Sections ?? and ?? Theorems 2.1.1 and 2.1.3 are proved for a rather abstract alphabet  $\Sigma$  (specifically, for any Hausdorff topological space  $\Sigma$ ). The combinatorial methods, unlike all other approaches for deriving large deviations principles, are based upon point estimates. For example they bound the probabilities associated with each possible outcome of the empirical measure (also denoted as *type*, see the definition below). This turns out to be very useful for some statistical applications (an example is presented in Section 2.6).

### 2.1.1 The method of types and Sanov's theorem

Let  $X_1, X_2, \dots, X_n$  be a sequence of random variables (r.v.) which are independent, identically distributed (i.i.d.) according to the law  $\mu \in M_1(\Sigma)$ . Throughout,  $M_1(\Sigma)$  denotes the space of all probability measures (laws) on the alphabet  $\Sigma$ . Typically we shall specify the topology with which  $M_1(\Sigma)$  is equipped. However, recall that here  $\Sigma$  is a finite set, so  $M_1(\Sigma)$  is identified with the standard probability simplex in  $\mathbb{R}^{|\Sigma|}$ , i.e., the set of all real valued vectors with  $|\Sigma|$  non-negative components which sum to 1, and open sets in  $M_1(\Sigma)$  are induced by open sets in  $\mathbb{R}^{|\Sigma|}$ .

Let  $\Sigma_\mu$  denote the support of the law  $\mu$ , i.e., the set  $\Sigma_\mu \triangleq \{a_i : \mu(a_i) > 0\}$ . In general,  $\Sigma_\mu$  may be a strict subset of  $\Sigma$  (i.e.,  $\mu$  may assign zero probability to some of the symbols in  $\Sigma$ ). However, when considering one underlying measure  $\mu$  we may without loss of generality (w.l.o.g.) reduce  $\Sigma$  to  $\Sigma_\mu$  by ignoring all symbols which appear with zero probability. This indeed is implicitly assumed throughout this Section (in Section 2.6 below we encounter a scenario where we have to keep track of various support sets of the form of  $\Sigma_\mu$ ).

**Definition 2.1.1** The type  $L_n^{\mathbf{x}}$  of a sequence  $\mathbf{x} \triangleq \{x_1, \dots, x_n\}$  is the empirical measure (law) associated with this sequence. Specifically,  $L_n^{\mathbf{x}} = \{L_n^{\mathbf{x}}(a_1), \dots, L_n^{\mathbf{x}}(a_{|\Sigma|})\}$  is an element of  $M_1(\Sigma)$ , where

$$L_n^{\mathbf{x}}(a_i) = \frac{1}{n} \sum_{j=1}^n 1_{x_j=a_i} \triangleq \frac{1}{n} N(a_i|\mathbf{x}), \quad i = 1, \dots, |\Sigma| \quad (2.1.1)$$

Let  $\mathcal{L}_n$  denote the set of all possible types of sequences of length  $n$ . Thus,  $\mathcal{L}_n \triangleq \{\nu : \nu = L_n^{\mathbf{x}} \text{ for some } \mathbf{x}\} \subset \mathbb{R}^{|\Sigma|}$ . Note that the random empirical measure  $L_n^{\mathbf{X}}$  associated with the sequence  $\mathbf{X} \triangleq \{X_1, \dots, X_n\}$  must be an element of the set  $\mathcal{L}_n$ .

The usefulness of the notion of types for finite alphabets as well as the reason why it is not readily extended to more general alphabets are due to the following ‘volume’ and ‘approximation distance’ estimates:

**Lemma 2.1.1** (a).  $|\mathcal{L}_n| \leq (n+1)^{|\Sigma|}$

(b). For any probability vector  $\nu \in M_1(\Sigma)$

$$d_V(\nu, \mathcal{L}_n) \triangleq \inf_{\nu' \in \mathcal{L}_n} d_V(\nu, \nu') \leq \frac{|\Sigma|}{2n}, \quad (2.1.2)$$

where  $d_V(\nu, \nu') \triangleq \sup_{A \subset \Sigma} [\nu(A) - \nu'(A)]$  is the variational distance between the measures  $\nu$  and  $\nu'$ .

**Proof:** Note that any component of the vector  $L_n^{\mathbf{x}}$  belongs to the set  $\{\frac{0}{n}, \frac{1}{n}, \dots, \frac{n}{n}\}$  whose cardinality (size) is  $(n+1)$ . Part (a) of the lemma follows since the vector  $L_n^{\mathbf{x}}$  is specified by  $|\Sigma|$  such quantities each of which assumes at most  $(n+1)$  distinct values.

To prove part (b) observe that  $\mathcal{L}_n$  indeed contains all probability vectors which are composed of  $|\Sigma|$  elements from the set  $\{\frac{0}{n}, \frac{1}{n}, \dots, \frac{n}{n}\}$ . Thus, for any vector  $\nu \in M_1(\Sigma)$  there exists a vector  $\nu' \in \mathcal{L}_n$  with  $|\nu(a_i) - \nu'(a_i)| \leq \frac{1}{n}$  for  $i = 1, \dots, |\Sigma|$ . The bound of (2.1.2) now follows since for any discrete alphabet

$$d_V(\nu, \nu') \leq \frac{1}{2} \sum_i |\nu(a_i) - \nu'(a_i)|.$$

□

**Remarks:**

1. Actually, since  $L_n^{\mathbf{x}}$  is a probability vector, at most  $|\Sigma| - 1$  such quantities should be specified and so  $|\mathcal{L}_n| \leq (n+1)^{|\Sigma|-1}$ .
2. Lemma 2.1.1 states that the volume of  $\mathcal{L}_n$ , the support of the random empirical measures  $L_n^{\mathbf{x}}$ , grows polynomially in  $n$  and further that for large enough  $n$ , the set  $\mathcal{L}_n$  approximates uniformly and arbitrarily well (in the sense of variational distance) any measure in  $M_1(\Sigma)$ . Both properties are invalid when  $|\Sigma| = \infty$ .

**Definition 2.1.2** *The type (composition) class  $T(\nu)$  of a probability law  $\nu \in \mathcal{L}_n$  is the set*  
 $T(\nu) = \{\mathbf{x} \in \Sigma^n : L_n^{\mathbf{x}} = \nu\}.$

Outcomes  $\mathbf{x}$  in the same type class are equally likely when  $X_1, \dots, X_n$  are *i.i.d.* random variables. This is part of the first among the three lemmas in the sequel which yield a strong prelude to large deviations results by estimating the exponential growth of each type class and the precise probability of any specific outcome from a given type class. The following definitions are useful for that purpose:

**Definitions:**

(1). The entropy of a probability vector  $\nu$  is  $H(\nu) \triangleq -\sum_{i=1}^{|\Sigma|} \nu(a_i) \log \nu(a_i)$ , where  $0 \log 0 = 0$  throughout.

(2). The (relative) entropy of a probability vector  $\nu$  relative to  $\mu$  is

$$H(\nu|\mu) \triangleq \sum_{a_i \in \Sigma_\nu} \nu(a_i) \log \frac{\nu(a_i)}{\mu(a_i)},$$

when  $\Sigma_\nu \subseteq \Sigma_\mu$  (in which case  $H(\nu|\mu) < \infty$ ). Otherwise  $H(\nu|\mu) \triangleq \infty$  (whenever  $\nu(a_i) > 0$  while  $\mu(a_i) = 0$  for some  $a_i \in \Sigma$ ).

**Remark:** Note that  $H(\cdot|\mu)$  is a rate function. Indeed, the set  $M_1(\Sigma_\mu)$  is a closed subset of  $M_1(\Sigma)$  in which  $H(\cdot|\mu)$  is a continuous function (note that the function  $-x \log x$  is continuous on  $[0,1]$ ), while  $H(\cdot|\mu) = \infty$  outside  $M_1(\Sigma_\mu)$ . Finally, the non-negativity of  $H(\cdot|\mu)$  follows by Jensen's inequality.

**Lemma 2.1.2** *If  $\mathbf{x} \in T(\nu)$  for  $\nu \in \mathcal{L}_n$ , then*

$$\text{Prob}_\mu(\{X_1, \dots, X_n\} = \mathbf{x}) = e^{-n[H(\nu) + H(\nu|\mu)]} \quad (2.1.3)$$

where  $\text{Prob}_\mu$  corresponds to the probability law  $\mu^{Z+}$  associated with an infinite sequence of random variables  $\{X_j\}$  which are independent and identically distributed according to  $\mu \in M_1(\Sigma)$ .

**Proof:** Clearly (2.1.3) holds when  $H(\nu|\mu) = \infty$  as the random empirical measure  $L_n^{\mathbf{x}}$  concentrates on types  $\nu \in \mathcal{L}_n$  for which  $\Sigma_\nu \subseteq \Sigma_\mu$  (i.e.,  $H(\nu|\mu) < \infty$ ). We may thus assume that  $H(\nu|\mu) < \infty$  and  $\Sigma_\nu \subseteq \Sigma_\mu$ .

Since  $L_n^{\mathbf{x}} = \nu$  and  $H(\nu) + H(\nu|\mu) = -\sum_{i=1}^{|\Sigma|} \nu(a_i) \log \mu(a_i)$  it follows that

$$\text{Prob}_\mu(\{X_1, \dots, X_n\} = \mathbf{x}) = \prod_{i=1}^{|\Sigma|} \mu(a_i)^{N(a_i|\mathbf{x})} = \prod_{i=1}^{|\Sigma|} \mu(a_i)^{n\nu(a_i)} = e^{-n[H(\nu) + H(\nu|\mu)]} \quad (2.1.4)$$

□

In particular, since  $H(\mu|\mu) = 0$ , when  $\mu \in \mathcal{L}_n$  and  $\mathbf{x} \in T(\mu)$  then

$$\text{Prob}_\mu(\{X_1, \dots, X_n\} = \mathbf{x}) = e^{-nH(\mu)}. \quad (2.1.5)$$

**Lemma 2.1.3** *For any  $\nu \in \mathcal{L}_n$ ,*  $(n+1)^{-|\Sigma|} e^{nH(\nu)} \leq |T(\nu)| \leq e^{nH(\nu)}$

**Remark:** Since  $|T(\nu)| = \binom{n}{n\nu(a_1), \dots, n\nu(a_{|\Sigma|})}$ , one might get a good estimate of  $|T(\nu)|$  by applying Stirling's approximation (see [12], pg 48). We take here a somewhat different route, of a more information theoretic flavor.

**Proof:** Any type class has probability at most one and all its members are of equal probability. Therefore, for any  $\nu \in \mathcal{L}_n$ , by (2.1.5),

$$1 \geq \text{Prob}_\nu(L_n^X = \nu) = \text{Prob}_\nu(X \in T(\nu)) = e^{-nH(\nu)} |T(\nu)| \quad (2.1.6)$$

and the upper bound on  $|T(\nu)|$  follows. The lower bound follows from the inequality (which is proved below)

$$\text{Prob}_\nu(L_n^X = \nu) \geq \text{Prob}_\nu(L_n^X = \nu'), \quad \forall \nu', \nu \in \mathcal{L}_n, \quad (2.1.7)$$

as thus

$$1 = \sum_{\nu' \in \mathcal{L}_n} \text{Prob}_\nu(L_n^X = \nu') \leq |\mathcal{L}_n| \text{Prob}_\nu(L_n^X = \nu) = |\mathcal{L}_n| e^{-nH(\nu)} |T(\nu)|,$$

while  $|\mathcal{L}_n| \leq (n+1)^{|\Sigma|}$  by Lemma 2.1.1.

It suffices to prove the inequality (2.1.7) for  $\nu' \in \mathcal{L}_n$  such that  $\Sigma_{\nu'} \subset \Sigma_\nu$  (as otherwise  $\text{Prob}_\nu(L_n^X = \nu') = 0$ ). Thus, without loss of generality one may assume that  $\Sigma = \Sigma_\nu$ . Now, consider the ratio

$$\begin{aligned} \frac{\text{Prob}_\nu(L_n^X = \nu)}{\text{Prob}_\nu(L_n^X = \nu')} &= \frac{\binom{n}{n\nu(a_1), \dots, n\nu(a_{|\Sigma|})}}{\binom{n}{n\nu'(a_1), \dots, n\nu'(a_{|\Sigma|})}} \frac{\prod_{i=1}^{|\Sigma|} \nu(a_i)^{n\nu(a_i)}}{\prod_{i=1}^{|\Sigma|} \nu(a_i)^{n\nu'(a_i)}} = \\ &= \prod_{i=1}^{|\Sigma|} \frac{(n\nu'(a_i))!}{(n\nu(a_i))!} \nu(a_i)^{n\nu(a_i) - n\nu'(a_i)}. \end{aligned} \quad (2.1.8)$$

This last expression is a product of terms of the form  $\frac{m!}{l!} \left(\frac{l}{n}\right)^{l-m}$ . By induction  $\frac{m!}{l!} \geq l^{(m-l)}$  for any  $m, l \in \mathbb{Z}^+$ , and thus (2.1.8) yields

$$\frac{\text{Prob}_\nu(L_n^X = \nu)}{\text{Prob}_\nu(L_n^X = \nu')} \geq \prod_{i=1}^{|\Sigma|} n^{n\nu'(a_i) - n\nu(a_i)} = n^{\sum_{i=1}^{|\Sigma|} \nu'(a_i) - \sum_{i=1}^{|\Sigma|} \nu(a_i)} = 1 \quad (2.1.9)$$

□



**Lemma 2.1.4 (Large deviations probabilities)** For any  $\nu \in \mathcal{L}_n$

$$(n+1)^{-|\Sigma|} e^{-nH(\nu|\mu)} \leq \text{Prob}_\mu(L_n^{\mathbf{X}} = \nu) \leq e^{-nH(\nu|\mu)} \quad (2.1.10)$$

**Proof:** By Lemma 2.1.2

$$\begin{aligned} \text{Prob}_\mu(L_n^{\mathbf{X}} = \nu) &= |T(\nu)| \text{Prob}_\mu(\{X_1, \dots, X_n\} = \mathbf{x} > L_n^{\mathbf{x}} = \nu) \\ &= |T(\nu)| e^{-n[H(\nu) + H(\nu|\mu)]} \end{aligned} \quad (2.1.11)$$

The proof is now completed by applying Lemma 2.1.3. □

Combining Lemma 2.1.1 and Lemma 2.1.4 we prove Sanov's theorem for the finite alphabet  $\Sigma$ , our first large deviations principle. Later on, in Sections ?? and ?? we shall extend this result to an appropriate class of topological spaces  $M_1(\Sigma)$  covering in particular the important case of  $\Sigma = \mathbb{R}^d$ .

Sanov's theorem deals with the sequence of laws governing the random empirical measure  $L_n^{\mathbf{X}}$  which is a *random element* of the probability simplex  $M_1(\Sigma)$  when the underlying independent random variables  $X_1, \dots, X_n$  are identically distributed according to the law  $\mu$ .

**Theorem 2.1.1 (Sanov's theorem)** For any set  $\Gamma$  of probability vectors in  $M_1(\Sigma) \subset \mathbb{R}^{|\Sigma|}$

$$-\inf_{\nu \in \Gamma^\circ} H(\nu|\mu) \leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \text{Prob}_\mu(L_n^{\mathbf{X}} \in \Gamma) \leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log \text{Prob}_\mu(L_n^{\mathbf{X}} \in \Gamma) \leq -\inf_{\nu \in \Gamma} H(\nu|\mu), \quad (2.1.12)$$

where  $\Gamma^\circ$  is the interior of  $\Gamma$  (as a set in  $M_1(\Sigma) \subset \mathbb{R}^{|\Sigma|}$ ).

**Remark:** Comparing (2.1.12) and the definition (1.1.13) we conclude that Sanov's theorem states that the family of laws  $\text{Prob}_\mu(L_n^{\mathbf{X}} \in \cdot)$  satisfies a large deviations principle with the rate function  $H(\cdot|\mu)$ . Further, in the particular case of finite alphabet  $\Sigma$  covered here the upper bound holds for any set  $\Gamma$  (i.e., there is no need for a closure operation). For few other improvements over (2.1.12) which are specific to this case see exercises 2.1.1, 2.1.2 and 2.1.3. Note however that there are closed sets for which the upper and lower bounds of (2.1.12) are far apart (see exercise 2.1.5). Moreover, there are closed sets for which the limit of  $\frac{1}{n} \log \text{Prob}_\mu(L_n^{\mathbf{X}} \in \Gamma)$  does not exist at all (see exercise 2.1.4).

**Proof:** We begin by deducing from Lemma 2.1.4 upper and lower bounds for any finite  $n$ . The upper bound

$$\begin{aligned} \text{Prob}_\mu(L_n^{\mathbf{X}} \in \Gamma) &= \sum_{\nu \in \Gamma \cap \mathcal{L}_n} \text{Prob}_\mu(L_n^{\mathbf{X}} = \nu) \leq \sum_{\nu \in \Gamma \cap \mathcal{L}_n} e^{-nH(\nu|\mu)} \leq |\Gamma \cap \mathcal{L}_n| e^{-n \inf_{\nu \in \Gamma \cap \mathcal{L}_n} H(\nu|\mu)} \\ &\leq (n+1)^{|\Sigma|} e^{-n \inf_{\nu \in \Gamma \cap \mathcal{L}_n} H(\nu|\mu)} . \end{aligned} \quad (2.1.13)$$

follows by the union of events bound and the upper bound of Lemma 2.1.4. The accompanying lower bound is

$$\begin{aligned} \text{Prob}_\mu(L_n^{\mathbf{X}} \in \Gamma) &= \sum_{\nu \in \Gamma \cap \mathcal{L}_n} \text{Prob}_\mu(L_n^{\mathbf{X}} = \nu) \geq \sum_{\nu \in \Gamma \cap \mathcal{L}_n} (n+1)^{-|\Sigma|} e^{-nH(\nu|\mu)} \geq \\ &\geq (n+1)^{-|\Sigma|} e^{-n \inf_{\nu \in \Gamma \cap \mathcal{L}_n} H(\nu|\mu)} . \end{aligned} \quad (2.1.14)$$

As  $\lim_{n \rightarrow \infty} \frac{1}{n} \log(n+1)^{|\Sigma|} = 0$  the normalized logarithmic limits of (2.1.13) and (2.1.14) yield

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log \text{Prob}_\mu(L_n^{\mathbf{X}} \in \Gamma) = - \liminf_{n \rightarrow \infty} \left\{ \inf_{\nu \in \Gamma \cap \mathcal{L}_n} H(\nu|\mu) \right\} \quad (2.1.15)$$

and

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log \text{Prob}_\mu(L_n^{\mathbf{X}} \in \Gamma) = - \limsup_{n \rightarrow \infty} \left\{ \inf_{\nu \in \Gamma \cap \mathcal{L}_n} H(\nu|\mu) \right\} . \quad (2.1.16)$$

The upper bound of (2.1.12) follows since  $\Gamma \cap \mathcal{L}_n \subset \Gamma$  for any  $n$ .

Turning now to the lower bound of (2.1.12), fix an arbitrary point  $\nu \in \Gamma^\circ$ . The set  $B_{\nu,\delta}^V \triangleq \{\nu' : d_V(\nu, \nu') < \delta\}$  must be contained in  $\Gamma$  for all  $\delta > 0$  sufficiently small (as  $\nu$  is in the interior of  $\Gamma$ ). Thus, by part (b) of Lemma 2.1.1 (see (2.1.2)) there exists a sequence  $\nu_n \in \Gamma \cap \mathcal{L}_n$  such that  $\nu_n \rightarrow \nu$  as  $n \rightarrow \infty$ . As already observed,  $\Sigma$  may be identified with  $\Sigma_\mu$  without loss of generality. Then  $H(\cdot|\mu)$  is a continuous function and therefore,

$$- \limsup_{n \rightarrow \infty} \inf_{\nu' \in \Gamma \cap \mathcal{L}_n} H(\nu'|\mu) \geq - \lim_{n \rightarrow \infty} H(\nu_n|\mu) = -H(\nu|\mu) . \quad (2.1.17)$$

The lower bound of (2.1.12) follows by taking the infimum over  $\nu \in \Gamma^\circ$  in (2.1.17).  $\square$

## Exercises:

### 2.1.1 Prove that for any open set $\Gamma$

$$- \lim_{n \rightarrow \infty} \left\{ \inf_{\nu \in \Gamma \cap \mathcal{L}_n} H(\nu|\mu) \right\} = \lim_{n \rightarrow \infty} \frac{1}{n} \log \text{Prob}_\mu(L_n^{\mathbf{X}} \in \Gamma) = - \inf_{\nu \in \Gamma} H(\nu|\mu) \triangleq -I_\Gamma . \quad (2.1.18)$$

2.1.2 (a). Extend the conclusions of exercise 2.1.1 to any set  $\Gamma \subseteq M_1(\Sigma_\mu)$  which is contained in the closure of its interior.

(b). Prove that for any such set  $I_\Gamma < \infty$  and  $I_\Gamma = H(\nu^*|\mu)$  for some  $\nu^* \in \overline{\Gamma^\circ}$ .

Hint: Use the continuity of  $H(\cdot|\mu)$  on  $M_1(\Sigma_\mu)$  and the compactness of this set.

2.1.3 Assume that  $\Gamma$  is a convex subset of  $M_1(\Sigma_\mu)$  of non-empty interior. Prove that all the conclusions of exercise 2.1.2 apply. Moreover, prove that here the point  $\nu^* \in \overline{\Gamma^\circ}$  for which  $I_\Gamma = H(\nu^*|\mu)$  is unique.

Hint: Show that for any  $\nu \in \Gamma$  and any  $\nu' \in \Gamma^\circ$  the line segment connecting  $\nu$  and  $\nu'$  is part of  $\Gamma^\circ$ . Deduce that  $\Gamma \subset \overline{\Gamma^\circ}$  which is a closed convex set. Then prove that  $H(\cdot|\mu)$  is a strictly convex function in  $M_1(\Sigma_\mu)$ .

2.1.4 Find a closed set  $\Gamma$  for which both of the limits in (2.1.18) do not exist.

Hint: Any set  $\Gamma = \{\nu\}$  consisting of one probability vector  $\nu \in M_1(\Sigma_\mu)$  which also belongs to  $\mathcal{L}_n$  (for some  $n$ ) shall do.

2.1.5 Find a closed set  $\Gamma$  such that  $\Gamma = \overline{\Gamma^\circ}$  and  $\inf_{\nu \in \Gamma} H(\nu|\mu) < \infty$  while  $\inf_{\nu \in \Gamma^\circ} H(\nu|\mu) = \infty$ .

Hint: For this construction you need  $\Sigma_\mu \neq \Sigma$ . Try  $\Gamma = \{\nu : \nu(a_2) + \sqrt{\nu(a_1)} \geq 1\}$  and  $\mu(a_1) = 0$  (while  $|\Sigma| = 3$ ).

## 2.1.2 Cramer's theorem for finite alphabets in $\mathbb{R}^1$

As an application of Sanov's theorem, we bring a proof of a version of Cramer's theorem about the large deviations of the empirical mean of i.i.d random variables. For that purpose we further assume throughout this section that  $\Sigma = \Sigma_\mu$  is a finite subset of  $\mathbb{R}^1$  and let  $\hat{S}_n \triangleq \frac{1}{n} \sum_{j=1}^n X_j$  denote the resulting sequence of empirical means (where  $X_j \in \Sigma$  as in Section 2.1.1 above). Cramer's theorem deals with the large deviations principle satisfied by the family of laws governing the real valued random variables  $\hat{S}_n$ . Sections 2.2 and 2.3 are devoted to successive generalizations of this result to  $\Sigma = \mathbb{R}^1$  (Section 2.2) and to weakly dependent random vectors in  $\mathbb{R}^d$  (Section 2.3).

Note that in the case considered here ( $|\Sigma| < \infty$ ), the random variables  $\hat{S}_n$  assume values in the compact interval  $K \triangleq [\min_{i=1}^{|\Sigma|} \{a_i\}, \max_{i=1}^{|\Sigma|} \{a_i\}]$ . Moreover,

$\hat{S}_n = \sum_{i=1}^{|\Sigma|} L_n^X(a_i) a_i = \langle L_n^X, \mathbf{a} \rangle$  where  $\mathbf{a} \triangleq (a_1, \dots, a_{|\Sigma|})$ . Therefore,

$$\hat{S}_n \in A \Leftrightarrow L_n^X \in \{\nu : \langle \nu, \mathbf{a} \rangle \in A\} \triangleq \Gamma \quad (2.1.19)$$

for any  $A \subset K$  and any integer  $n$ . Thus, the following version of Cramer's theorem is a direct consequence of Theorem 2.1.1.

**Theorem 2.1.2 (Cramer's theorem for finite subsets of  $\mathbb{R}^1$ )**

(a). For any set  $A \subset K$

$$-\inf_{x \in A^\circ} I(x) \leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \text{Prob}_\mu(\hat{S}_n \in A) \leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log \text{Prob}_\mu(\hat{S}_n \in A) \leq -\inf_{x \in A} I(x), \quad (2.1.20)$$

where  $A^\circ$  is the interior of  $A$  (as a set in  $\mathbb{R}^1$ ) and  $I(x) \triangleq \inf_{\langle \nu, \mathbf{a} \rangle = x} H(\nu | \mu)$ .

(b). The continuous rate function  $I(x)$  (for  $x \in K$ ) also satisfies

$$I(x) = \sup_{\lambda \in \mathbb{R}^1} \{\lambda x - \Lambda(\lambda)\}, \quad (2.1.21)$$

where

$$\Lambda(\lambda) = \log \sum_{i=1}^{|\Sigma|} \mu(a_i) e^{\lambda a_i}.$$

**Remark:** Since the rate function is continuous (on  $K$ ) it follows from (2.1.20) that whenever  $A \subset \overline{A^\circ} \subset K$

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \text{Prob}_\mu(\hat{S}_n \in A) = -\inf_{x \in A} I(x).$$

**Proof:**

(a). The bounds of (2.1.20) are simply the bounds of (2.1.12) for the set  $\Gamma$  defined in (2.1.19) (note that  $\{\nu : \langle \nu, \mathbf{a} \rangle \in A^\circ\} \subset \Gamma^\circ$ ).

(b). Observe that by Jensen's inequality

$$\Lambda(\lambda) = \log \sum_{i=1}^{|\Sigma|} \mu(a_i) e^{\lambda a_i} \geq \sum_{i=1}^{|\Sigma|} \nu(a_i) \log \frac{\mu(a_i) e^{\lambda a_i}}{\nu(a_i)} = \lambda \langle \nu, \mathbf{a} \rangle - H(\nu | \mu), \quad (2.1.22)$$

for any  $\nu \in M_1(\Sigma)$  and any  $\lambda \in \mathbb{R}^1$ , with equality for  $\nu_\lambda(a_i) \triangleq \mu(a_i) e^{\lambda a_i - \Lambda(\lambda)}$ .

The function  $\Lambda(\lambda)$  is differentiable and strictly convex, implying that  $\Lambda'(\lambda)$  is strictly increasing. Moreover,  $K = [\inf_\lambda \Lambda'(\lambda), \sup_\lambda \Lambda'(\lambda)]$ . Therefore, for any  $x \in K$  there exists a unique  $\lambda^{(x)} \in \mathbb{R}^1$

such that  $x - \Lambda'(\lambda^{(x)}) = 0$ . As  $\langle \nu_{\lambda^{(x)}} \rangle = \Lambda'(\lambda^{(x)}) = x$ , the application of (2.1.22) for  $\lambda^{(x)}$  yields

$$I(x) = \inf_{\langle \nu, \mathbf{a} \rangle = x} H(\nu|\mu) = H(\nu_{\lambda^{(x)}}|\mu) = \lambda^{(x)}x - \Lambda(\lambda^{(x)}) = \sup_{\lambda \in \mathbb{R}^1} \{\lambda x - \Lambda(\lambda)\}, \quad (2.1.23)$$

i.e., (2.1.21) holds for any  $x \in K^\circ$ . The continuity of  $I(x)$  (for  $x \in K$ ), which is a direct consequence of the continuity of  $H(\cdot|\mu)$ , implies that (2.1.21) holds also at the boundaries of  $K$ .  $\square$

It is interesting to note that by deriving Cramer's theorem from Sanov's theorem, we have followed a pattern which will be useful in the sequel and which is referred to as the 'contraction principle' (cf Section ??). Indeed, based on large deviations result in the "big" space  $(M_1(\Sigma))$ , we obtained a large deviations result on a "smaller" space which is obtained from the original space by a continuous map ("contraction"). In the particular case in hand, however, an alternative proof of the finite alphabet Sanov's theorem follows directly from Cramer's theorem in  $\mathbb{R}^{|\Sigma|}$  and is presented in section 2.3.

#### Exercises:

**2.1.6** Construct an example for which the limit of  $\frac{1}{n} \log \text{Prob}_\mu(\hat{S}_n = x)$  as  $n \rightarrow \infty$  does not exist. **Hint:** Note that for  $|\Sigma| = 2$  the empirical mean  $(\hat{S}_n)$  uniquely determines the empirical measure  $(L_n^X)$ . Rely on this observation and exercise 2.1.4 to construct the desired example.

**2.1.7 (a).** Prove that  $I(x) = 0$  if and only if  $x = \bar{x} \triangleq E_\mu(X_1)$ . Explain why this should have been anticipated in view of the weak law of large numbers (which states that  $\hat{S}_n \rightarrow \bar{x}$  in probability, as  $n \rightarrow \infty$ ).

(b). Check that  $H(\nu|\mu) = 0$  if and only if  $\nu = \mu$  and interpret this result.

**2.1.8** Guess the value of  $\lim_{n \rightarrow \infty} \text{Prob}_\mu(X_1 = a_i | \hat{S}_n \geq q)$  for  $\bar{x} < q \in K^\circ$  and try to justify your guess (at least heuristically).

**2.1.9** Extend Theorem 2.1.2 to  $\Sigma$  which is a finite subset of  $\mathbb{R}^d$ . In particular, determine the shape of the set  $K$  and find the appropriate extension of the formula (2.1.21).

### 2.1.3 Large deviations for sampling without replacements

The method of types is useful also for studying the large deviations of the empirical measure of a sequence of dependent random variables. For example, consider the following set up of sampling

without replacements which is often encountered in statistical problems. Out of an initial deterministic pool of  $M$  items  $\mathbf{x} \triangleq (x_1, \dots, x_M)$ , an  $n$ -tuple  $\mathbf{X} \triangleq (x_{i_1}, x_{i_2}, \dots, x_{i_n})$  is sampled without replacement, namely,  $i_1 \neq i_2 \neq i_3 \neq \dots \neq i_n$  and all  $\binom{M}{n}$  choices of  $i_1 \neq i_2 \dots \neq i_n \in \{1, \dots, M\}$  are equally likely.

Suppose now that  $x_1^{(M)}, \dots, x_M^{(M)}$  are elements of the same finite set  $\Sigma = \{a_1, \dots, a_{|\Sigma|}\}$  and that as  $M \rightarrow \infty$  the deterministic relative frequency vectors  $L_M^{\mathbf{x}} = \{L_M^{\mathbf{x}}(a_1), \dots, L_M^{\mathbf{x}}(a_{|\Sigma|})\}$  converge (point wise) to  $\mu \in M_1(\Sigma)$ . Recall that

$$L_M^{\mathbf{x}}(a_i) = \frac{1}{M} \sum_{j=1}^M 1_{x_j^{(M)}=a_i} \triangleq \frac{1}{M} M(a_i | \mathbf{x}), \quad i = 1, \dots, |\Sigma|. \quad (2.1.24)$$

Suppose further that  $\mathbf{X}$  is a random vector obtained by the sampling without replacement of  $n$  out of  $M$  elements as described above. We investigate the large deviations principle governing the random empirical measure  $L_n^{\mathbf{X}}$  associated with the vector  $\mathbf{X}$ . Note that  $L_n^{\mathbf{X}}$  belongs to the set  $\mathcal{L}_n$  whose size grows polynomially in  $n$  (see Lemma 2.1.1). In particular we aim at deriving the analog of Theorem 2.1.1 when  $\lim_{M \rightarrow \infty} (\frac{n}{M}) = \beta \in (0, 1)$ . For that purpose define the following candidate rate function

$$I_{\beta, \mu}(\nu) \triangleq \begin{cases} H(\nu | \mu) + \frac{1-\beta}{\beta} H\left(\frac{\mu - \beta \nu}{1-\beta} \middle| \mu\right) & \text{if } \mu_i \geq \beta \nu_i \text{ for all } i \\ \infty & \text{otherwise} \end{cases} \quad (2.1.25)$$

By basic combinatorics and an application of Lemma 2.1.3 we obtain the following estimates of large deviations probabilities for  $L_n^{\mathbf{X}}$ .

**Lemma 2.1.5** *For any probability vector  $\nu \in \mathcal{L}_n$*

(a). *If  $I_{\frac{n}{M}, L_M^{\mathbf{x}}}(\nu) < \infty$  then*

$$\left| \frac{1}{n} \log \text{Prob}(L_n^{\mathbf{X}} = \nu) + I_{\frac{n}{M}, L_M^{\mathbf{x}}}(\nu) \right| \leq 2(|\Sigma| + 1) \left( \frac{\log(M+1)}{n} \right) \quad (2.1.26)$$

(b). *If  $I_{\frac{n}{M}, L_M^{\mathbf{x}}}(\nu) = \infty$  then*

$$\text{Prob}(L_n^{\mathbf{X}} = \nu) = 0 \quad (2.1.27)$$

**Proof:** (a). In the sampling without replacement procedure the probability of the event  $\{L_n^{\mathbf{X}} = \nu\}$  for  $\nu \in \mathcal{L}_n$  is exactly the number of  $n$ -tuples  $i_1 \neq i_2 \dots \neq i_n$  resulting with type  $\nu$  divided by the

overall number of  $n$ -tuples, that is

$$\text{Prob}(L_n^{\mathbf{X}} = \nu) = \frac{\prod_{i=1}^{|\Sigma|} \binom{M L_M^{\mathbf{X}}(a_i)}{n \nu(a_i)}}{\binom{M}{n}}. \quad (2.1.28)$$

An application of Lemma 2.1.3 for  $|\Sigma| = 2$  (where  $|T(\nu)| = \binom{n}{k}$  when  $\nu(a_1) = \frac{k}{n}$ ,  $\nu(a_2) = 1 - \frac{k}{n}$ ), results with the following estimate

$$\max_{0 \leq k \leq n} \left| \log \binom{n}{k} - n H\left(\frac{k}{n}\right) \right| \leq 2 \log(n+1), \quad (2.1.29)$$

where

$$H(p) \triangleq -p \log p - (1-p) \log(1-p).$$

Alternatively, (2.1.29) follows by Stirling's formula (see [12], pg 48).

Combining the exact expression (2.1.28) and the bound of (2.1.29) results with

$$\left| \frac{1}{n} \log \text{Prob}(L_n^{\mathbf{X}} = \nu) - \sum_{i=1}^{|\Sigma|} \frac{M L_M^{\mathbf{X}}(a_i)}{n} H\left(\frac{n \nu(a_i)}{M L_M^{\mathbf{X}}(a_i)}\right) + \frac{M}{n} H\left(\frac{n}{M}\right) \right| \leq 2(|\Sigma|+1) \left( \frac{\log(M+1)}{n} \right) \quad (2.1.30)$$

The inequality (2.1.26) follows when rearranging the left side of (2.1.30).

(b). Note that  $I_{\frac{n}{M}, L_M^{\mathbf{X}}}(\nu) = \infty$  if and only if  $n \nu(a_i) > M L_M^{\mathbf{X}}(a_i) = M(a_i|\mathbf{x})$  for some  $a_i \in \Sigma$ . It is then impossible in sampling without replacement to have  $L_n^{\mathbf{X}} = \nu$ , as  $n L_n^{\mathbf{X}}(a_i) = N(a_i|\mathbf{X}) \leq M(a_i|\mathbf{x})$  for any  $a_i \in \Sigma$ .  $\square$

Following the proof of Theorem 2.1.1 the above point estimates result with the analogs of (2.1.15) and (2.1.16), namely:

**Corollary 2.1.1**

$$\limsup_{M \rightarrow \infty} \frac{1}{n} \log \text{Prob}(L_n^{\mathbf{X}} \in \Gamma) = - \liminf_{M \rightarrow \infty} \left\{ \inf_{\nu \in \Gamma \cap \mathcal{L}_n} I_{\frac{n}{M}, L_M^{\mathbf{X}}}(\nu) \right\} \quad (2.1.31)$$

and

$$\liminf_{M \rightarrow \infty} \frac{1}{n} \log \text{Prob}(L_n^{\mathbf{X}} \in \Gamma) = - \limsup_{M \rightarrow \infty} \left\{ \inf_{\nu \in \Gamma \cap \mathcal{L}_n} I_{\frac{n}{M}, L_M^{\mathbf{X}}}(\nu) \right\}. \quad (2.1.32)$$

The proof of Corollary 2.1.1 is left as exercise 2.1.10.

The following theorem is the desired analog of Theorem 2.1.1.

**Theorem 2.1.3** *The good rate function  $I_{\beta,\mu}(\nu)$  controls the large deviations principle related to the random empirical measures  $L_n^{\mathbf{X}}$  (as elements of  $M_1(\Sigma)$ ). Specifically, for any set  $\Gamma$  of probability vectors in  $M_1(\Sigma) \subset \mathbb{R}^{|\Sigma|}$ ,*

$$-\inf_{\nu \in \Gamma^o} I_{\beta,\mu}(\nu) \leq \liminf_{M \rightarrow \infty} \frac{1}{n} \log \text{Prob}(L_n^{\mathbf{X}} \in \Gamma) \leq \limsup_{M \rightarrow \infty} \frac{1}{n} \log \text{Prob}(L_n^{\mathbf{X}} \in \Gamma) \leq -\inf_{\nu \in \Gamma} I_{\beta,\mu}(\nu), \quad (2.1.33)$$

where  $\beta = \lim_{M \rightarrow \infty} (\frac{n}{M}) \in (0, 1)$  and  $L_M^{\mathbf{X}}$  converges to  $\mu$  as  $M \rightarrow \infty$ .

**Remark:** Note that the upper bound of (2.1.33) is weaker than the upper bound of (2.1.12). Also, see exercise 2.1.11 for examples of sets for which the above lower and upper bounds coincide.

**Proof:** As already noted,  $H(\cdot|\mu)$  is a rate function on  $M_1(\Sigma)$ , for any fixed  $\mu$  (namely, it is lower semicontinuous and non-negative). Moreover, the set  $\{\nu : (\mu - \beta\nu)/(1 - \beta) \in M_1(\Sigma)\}$  is a closed subset of  $M_1(\Sigma)$  and therefore,  $I_{\beta,\mu}(\cdot)$  is also a rate function for any fixed  $\beta \in (0, 1)$  and any fixed  $\mu \in M_1(\Sigma)$ . Since  $|\Sigma| < \infty$ , the probability simplex  $M_1(\Sigma)$  is a compact set and thus any rate function on  $M_1(\Sigma)$  is a good rate function.

The first step in deriving the bounds of (2.1.33) is to prove that  $I_{\beta,\mu}(\nu)$  is a lower semicontinuous function over  $(0, 1) \times M_1(\Sigma) \times M_1(\Sigma)$  (jointly in  $\beta$ ,  $\mu$  and  $\nu$ ) and is strictly continuous along sequences  $\{\beta_n, \mu_n, \nu_n\}$  where  $I_{\beta_n, \mu_n}(\nu_n) < \infty$ . For that purpose, fix a sequence  $\{\beta_n, \mu_n, \nu_n\}$  such that  $\beta_n \rightarrow \beta \in (0, 1)$ ,  $\nu_n \rightarrow \nu$  and  $\mu_n \rightarrow \mu$ . The lower semicontinuity is trivial along any subsequence  $\{n_k\}_{k=1}^{\infty}$  for which  $I_{\beta_{n_k}, \mu_{n_k}}(\nu_{n_k}) = \infty$ . Thus, without loss of generality assume that  $I_{\beta_n, \mu_n}(\nu_n) < \infty$  for all  $n$  large enough, namely that  $\mu_n(a_i) \geq \beta_n \nu_n(a_i)$  for any  $a_i \in \Sigma$ . Then, by rearranging (2.1.25)

$$I_{\beta_n, \mu_n}(\nu_n) = \frac{1}{\beta_n} H(\mu_n) - H(\nu_n) - \frac{1 - \beta_n}{\beta_n} H\left(\frac{\mu_n - \beta_n \nu_n}{1 - \beta_n}\right).$$

As  $\beta_n$  are bounded away from 0 and 1 for large  $n$ , it follows from the above expression that  $I_{\beta_n, \mu_n}(\nu_n) \rightarrow I_{\beta, \mu}(\nu)$  since the entropy function  $H(\cdot)$  is strictly continuous over  $M_1(\Sigma)$  (for  $|\Sigma| < \infty$ ).

We turn now to prove the upper bound of (2.1.33). We first deduce from (2.1.31) that for some



infinite subsequence  $M_k$ , there exist  $\nu_k \in \Gamma$  such that

$$\limsup_{M \rightarrow \infty} \frac{1}{n} \log \text{Prob}(L_n^X \in \Gamma) = - \lim_{k \rightarrow \infty} I_{\frac{n}{M_k}, L_{M_k}^X}(\nu_k) \triangleq I^* . \quad (2.1.34)$$

Moreover, the sequence  $\{\nu_k\}$  has a limit point  $\nu^*$  in the compact set  $\bar{\Gamma}$ . Passing to a convergent subsequence, the lower semicontinuity of  $I$  jointly in  $\beta$ ,  $\mu$  and  $\nu$  implies that

$$I^* \leq -I_{\beta, \mu}(\nu^*) \leq - \inf_{\nu \in \bar{\Gamma}} I_{\beta, \mu}(\nu) . \quad (2.1.35)$$

The proof of the upper bound is complete in view of (2.1.34) and (2.1.35).

In view of (2.1.32) it suffices for the lower bound of (2.1.33) to prove that

$$- \limsup_{M \rightarrow \infty} \left\{ \inf_{\nu' \in \Gamma \cap \mathcal{L}_n} I_{\frac{n}{M}, L_M^X}(\nu') \right\} \geq -I_{\beta, \mu}(\nu) \quad (2.1.36)$$

for any  $\nu \in \Gamma^\circ$  where  $I_{\beta, \mu}(\nu) < \infty$  (see (1.1.7) in Section 1.1 for this formulation of the large deviations lower bound). Fix such a point  $\nu$ , i.e., such that

$$\mu(a_i) \geq \beta \nu(a_i) \quad \forall a_i \in \Sigma . \quad (2.1.37)$$

In particular,  $\Sigma_\nu \subseteq \Sigma_\mu$  and as  $\nu \in \Gamma^\circ$  as well,  $\nu$  belongs to the relative interior of  $\Gamma$  in  $M_1(\Sigma_\mu)$ . Further, (2.1.36) follows from

$$- \limsup_{M \rightarrow \infty} \left\{ \inf_{\nu' \in \Gamma \cap M_1(\Sigma_\mu) \cap \mathcal{L}_n} I_{\frac{n}{M}, L_M^X}(\nu') \right\} \geq -I_{\beta, \mu}(\nu) .$$

Therefore, one may assume without loss of generality that  $\Sigma = \Sigma_\mu$ . Now, since  $\beta < 1$  there always exists a sequence of probability vectors  $\nu_k \in \Gamma^\circ$  which converges to  $\nu$  and for which *all* the inequalities in (2.1.37) are *strict*. As  $I_{\beta, \mu}(\cdot)$  is continuous along any such sequence, we deduce that it suffices to prove (2.1.36) for  $\nu \in \Gamma^\circ$  for which

$$\min_{a_i \in \Sigma} \{\mu(a_i) - \beta \nu(a_i)\} > 0 . \quad (2.1.38)$$

By the same argument as in the proof of Theorem 2.1.1 there exists a sequence  $\nu_n \in \Gamma \cap \mathcal{L}_n$  such that  $\nu_n \rightarrow \nu$  as  $n \rightarrow \infty$ . Now, since  $L_M^X \rightarrow \mu$  and  $\frac{n}{M} \rightarrow \beta$  the *strict* inequality (2.1.38) implies that for all  $M$  sufficiently large

$$\min_{a_i \in \Sigma} \{L_M^X(a_i) - \frac{n}{M} \nu_n(a_i)\} \geq 0$$

and then  $I_{\frac{n}{M}, L_M^X}(\nu_n) < \infty$  (note also that  $\frac{n}{M}$  is eventually bounded away from 0 and 1). Therefore, by the strict continuity of  $I$  along such sequences

$$I_{\beta, \mu}(\nu) = \lim_{M \rightarrow \infty} I_{\frac{n}{M}, L_M^X}(\nu_n) .$$

Since  $\nu_n \in \Gamma \cap \mathcal{L}_n$ ,

$$\limsup_{M \rightarrow \infty} \left\{ \inf_{\nu \in \Gamma \cap \mathcal{L}_n} I_{\frac{n}{M}, L_M^X}(\nu) \right\} \leq \lim_{M \rightarrow \infty} I_{\frac{n}{M}, L_M^X}(\nu_n)$$

and the inequality (2.1.36) now follows.  $\square$

### Exercises:

2.1.10 Prove Corollary 2.1.1.

2.1.11 Let  $I_\Gamma \triangleq \inf_{\nu \in \Gamma^\circ} I_{\beta, \mu}(\nu)$ . Prove that when  $I_\Gamma < \infty$  and  $\Gamma$  is included in the closure of its interior, then

$$I_\Gamma = \lim_{M \rightarrow \infty} \frac{1}{n} \log \text{Prob}(L_n^X \in \Gamma) .$$

Hint: See exercise 2.1.2 and use the continuity of  $I_{\beta, \mu}(\cdot)$  within its level sets.

2.1.12 Prove that the rate function  $I_{\beta, \mu}(\cdot)$  is a convex function.

2.1.13 Prove that the large deviations principle of Theorem 2.1.3 holds for  $\beta = 0$  with the good rate function  $I_{0, \mu}(\cdot) \triangleq H(\cdot | \mu)$ , provided that  $n \rightarrow \infty$ .

Hint: First prove that the left side of (2.1.30) goes to zero as long as  $n \rightarrow \infty$  (i.e., even when  $\left(\frac{\log(M+1)}{n}\right) \rightarrow \infty$ ). Then prove the lower semicontinuity of  $I_{\beta, \mu}(\nu)$  at  $\beta = 0$  and use it to derive the upper bound. Finally, it suffices to prove the lower bound when  $\Sigma_\nu \subseteq \Sigma_\mu$  and so the sequence  $\nu_n \in M_1(\Sigma_\mu) \cap \Gamma \cap \mathcal{L}_n$  will converge to  $\nu$  while eventually  $I_{\frac{n}{M}, L_M^X}(\nu_n) < \infty$ .

## 2.2 Cramer's Theorem in $\mathbb{R}^1$

Cramer's theorem about the large deviations associated with the empirical means of i.i.d. random variables is presented in Section 2.1.2 as an application of the method of types. However, the method of types is limited in its scope to finite alphabets. Moreover, it neither explains why a large deviations principle is satisfied nor predicts which rate functions one should expect in similar

situations. In this section, we start pursuing a different route aiming at proving the analog of Theorem 2.1.2 when the underlying alphabet  $\Sigma$  is  $\mathbb{R}^1$ . The approach outlined here is more amenable for generalizations and better illustrates the main ingredients involved in proving a typical large deviations principle.

Consider the empirical means  $\hat{S}_n \triangleq \frac{1}{n} \sum_{j=1}^n X_j$ , where the real valued random variables  $X_1, \dots, X_n, \dots$  are independent and identically distributed according to the marginal probability law  $\mu \in M_1(\mathbb{R})$ . Let  $\mu_n$  denote the law of  $\hat{S}_n$  and  $\bar{x} \triangleq E_\mu[X_1]$  denotes the true underlying mean. If we further assume that  $\bar{x} < \infty$  and  $E[|X_1 - \bar{x}|^2] < \infty$  then  $\hat{S}_n \xrightarrow[n \rightarrow \infty]{Prob} \bar{x}$  since

$$E[(\hat{S}_n - \bar{x})^2] = \frac{1}{n^2} \sum_{j=1}^n E[|X_j - \bar{x}|^2] = \frac{1}{n} E[|X_1 - \bar{x}|^2] \xrightarrow[n \rightarrow \infty]{} 0. \quad (2.2.39)$$

Actually, the finite variance condition is not needed for the convergence in probability but we will not care about that here. Because of (2.2.39),  $\mu_n(A) \xrightarrow[n \rightarrow \infty]{} 0$  for any set  $A$  such that  $\bar{x} \notin \bar{A}$ . Cramer's theorem characterizes the logarithmic rate of this convergence by the (rate) function

$$\Lambda^*(x) \triangleq \sup_{\lambda \in \mathbb{R}^1} [\lambda x - \Lambda(\lambda)], \quad (2.2.40)$$

where

$$\Lambda(\lambda) \triangleq \log M(\lambda) \triangleq \log E[e^{\lambda X_1}], \quad (2.2.41)$$

is also called the **logarithmic moment generating function**. Note that while  $\Lambda(\lambda) > -\infty$  it is possible to have  $\Lambda(\lambda) = \infty$ . Let  $\mathcal{D}_\Lambda \triangleq \{\lambda : \Lambda(\lambda) < \infty\}$  and  $\mathcal{D}_{\Lambda^*} \triangleq \{x : \Lambda^*(x) < \infty\}$ .

The rate function  $\Lambda^*(x)$  of (2.2.40) is called the *Legendre transform* of  $\Lambda(\lambda)$ . Some of the properties of this function (and of  $\Lambda(\lambda)$ ) which are useful when proving Cramer's theorem are summarized in the following lemma whose proof is deferred to the end of this section. The exact definition of Legendre transform and its properties for more general vector spaces are presented in Sections 2.3 (for  $\mathbb{R}^d$ ) and ?? (for a general class of metric spaces).

### Lemma 2.2.1

- (a).  $\Lambda$  is a convex function and  $\Lambda^*$  is a convex rate function.
- (b).  $\Lambda^*(\bar{x}) = 0$  and for  $x \in [\bar{x}, \infty)$ ,

$$\Lambda^*(x) = \sup_{\lambda \geq 0} [\lambda x - \Lambda(\lambda)], \quad (2.2.42)$$

is nondecreasing. Similarly, for  $x \in (-\infty, \bar{x}]$

$$\Lambda^*(x) = \sup_{\lambda \leq 0} [\lambda x - \Lambda(\lambda)] \quad (2.2.43)$$

and is nonincreasing.

(c). For any  $\eta \in \mathcal{D}_\Lambda^\circ$

$$\Lambda'(\eta) = \frac{1}{M(\eta)} E[X_1 e^{\eta X_1}] \quad (2.2.44)$$

and

$$\Lambda'(\eta) = y \implies \Lambda^*(y) = \eta y - \Lambda(\eta) \quad (2.2.45)$$

The following definition is needed for the precise statement of Cramer's theorem.

**Definition 2.2.1** *The logarithmic moment generating function  $\Lambda(\cdot)$  is **steep** if  $\liminf_{n \rightarrow \infty} |\Lambda'(\lambda_n)| = \infty$  for any sequence  $\lambda_n \in \mathcal{D}_\Lambda^\circ$  which converges to  $\lambda \in \mathcal{D}_\Lambda \setminus \mathcal{D}_\Lambda^\circ$ .*

A weak version of Cramer's theorem is proved below, establishing the large deviations principle in  $\mathbb{R}^1$  when  $\Lambda(\cdot)$  is a steep function and  $0 \in \mathcal{D}_\Lambda^\circ$ . For example, both of these conditions hold when  $\mathcal{D}_\Lambda = \mathbb{R}^1$ . It is possible to establish the large deviations principle in  $\mathbb{R}^1$  with no condition on  $\Lambda(\lambda)$  by adopting a more sophisticated proof based on sub-additivity and convexity arguments. This is pursued in much generality in Sections ?? and ?? where the exact consequences for the special case discussed here are given in exercise ??.

### Theorem 2.2.1 (Cramer)

(a). *The large deviations upper bound*

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(A) \leq - \inf_{x \in A} \Lambda^*(x), \quad (2.2.46)$$

*holds for any measurable  $A \subset \mathbb{R}^1$ .*

(b). *For any open set  $G \subset \mathbb{R}^1$*

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(G) \geq - \inf_{x \in G \cap \bar{\mathcal{F}}} \Lambda^*(x), \quad (2.2.47)$$

*where  $\mathcal{F} \triangleq \{x : x = \Lambda'(\lambda) \text{ for some } \lambda \in \mathcal{D}_\Lambda^\circ\}$ .*

(c). *If  $\mathcal{D}_{\Lambda^*} \subseteq \bar{\mathcal{F}}$  then the family of probability measures  $\mu_n$  satisfies the large deviations principle*

with the rate function  $\Lambda^*(x)$ . In particular,  $\mu_n$  satisfies the large deviations principle if  $\mathcal{D}_\Lambda = \mathbb{R}^1$ .  
(d). If  $0 \in \mathcal{D}_\Lambda^\circ$  and  $\Lambda$  is a steep function then the family of probability measures  $\mu_n$  satisfies the large deviations principle with the good rate function  $\Lambda^*(x)$ .

**Remark:** The Legendre transform  $\Lambda^*(x)$  is a natural candidate for the rate function under a very general situation where the upper bound (2.2.46) holds (for more on this issue, see Sections 2.3 and ??). For that reason, convexity plays an important role in large deviations as seen throughout Chapters ??-??.

**Proof:**

(a). If  $\bar{x} \in \bar{A}$  then the upper bound (2.2.46) trivially holds (as  $\Lambda^*(\bar{x}) = 0$ ). Otherwise, let  $(x_-, x_+)$  be the largest interval around  $\bar{x}$  such that  $(x_-, x_+) \cap \bar{A} = \emptyset$ . Since

$$\bar{A} \subseteq (-\infty, x_-] \cup [x_+, \infty) \triangleq A_- \cup A_+,$$

it suffices to prove that

$$\mu_n(A_+) \leq e^{-n\Lambda^*(x_+)} \quad (2.2.48)$$

and

$$\mu_n(A_-) \leq e^{-n\Lambda^*(x_-)} \quad (2.2.49)$$

in order to prove the upper bound. Indeed, then

$$\mu_n(A) \leq \mu_n(A_+) + \mu_n(A_-) \leq e^{-n\Lambda^*(x_+)} + e^{-n\Lambda^*(x_-)} \leq 2e^{-n \min\{\Lambda^*(x_-), \Lambda^*(x_+)\}}.$$

By part (b) of Lemma 2.2.1

$$\min\{\Lambda^*(x_-), \Lambda^*(x_+)\} = \inf_{x \in \bar{A}} \Lambda^*(x)$$

and the proof of the upper bound follows.

Returning now to prove (2.2.48) and (2.2.49), these are consequences of a special form of Markov's inequality. Clearly, for any  $\lambda \geq 0$ ,

$$\begin{aligned} \mu_n(A_+) &\leq \int_{A_+} e^{n\lambda(z-x_+)} d\mu_n(z) \leq e^{-n\lambda x_+} E(e^{\lambda n \hat{S}_n}) \\ &= e^{-n\lambda x_+} \prod_{i=1}^n E(e^{\lambda X_i}) = e^{-n[\lambda x_+ - \Lambda(\lambda)]} \end{aligned} \quad (2.2.50)$$

Therefore,

$$\mu_n(A_+) \leq e^{-\sup_{\lambda \geq 0} n[\lambda x_+ - \Lambda(\lambda)]}$$

and (2.2.48) follows by (2.2.42) as  $x_+ \in [\bar{x}, \infty)$ . The proof of (2.2.49) is similar.

(b). The lower bound of (2.2.47) trivially holds when  $G \cap \bar{\mathcal{F}}$  is an empty set. Moreover, since  $G$  is an open set it suffices to show that for any  $x \in \bar{\mathcal{F}}$ , and any open interval  $B_{x,\delta}$  of center at  $x$  (and width  $2\delta$ )

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(B_{x,\delta}) \geq -\Lambda^*(x). \quad (2.2.51)$$

Consider first points  $y \in \mathcal{F}$ . Fix such a point  $y = \Lambda'(\eta) \in \mathcal{F}$  and fix  $\delta > 0$ . Note that  $\eta \in \mathcal{D}_\Lambda^\circ$  and define the associated probability measure (see [12])

$$\tilde{\mu}(dx) = e^{\eta x - \Lambda(\eta)} \mu(dx).$$

Indeed  $\tilde{\mu}$  is a probability measure as  $\int_{\mathbb{R}^1} \tilde{\mu}(dx) = \frac{1}{M(\eta)} \int_{\mathbb{R}^1} e^{\eta x} \mu(dx) = 1$ . Define accordingly  $\tilde{\mu}_n$  to be the law governing  $\hat{S}_n$  when  $X_i$  are i.i.d. random variables of law  $\tilde{\mu}$ . Note that by (2.2.44)

$$E_{\tilde{\mu}}(X_1) \triangleq \frac{1}{M(\eta)} \int_{\mathbb{R}^1} x e^{\eta x} \mu(dx) = \Lambda'(\eta) = y.$$

Thus, the mean of  $X_i$  under  $\tilde{\mu}$  equals  $y$  and by the weak law of large numbers,

$$\lim_{n \rightarrow \infty} \tilde{\mu}_n(B_{y,\delta}) = 1. \quad (2.2.52)$$

Moreover,

$$\begin{aligned} \mu_n(B_{y,\delta}) &= \int_{|\hat{S}_n - y| < \delta} d\mu(x_1) \cdots d\mu(x_n) \geq e^{-n(\eta y + \delta|\eta|)} \int_{|\hat{S}_n - y| < \delta} e^{\eta \sum_{i=1}^n x_i} d\mu(x_1) \cdots d\mu(x_n) \\ &= \tilde{\mu}_n(B_{y,\delta}) e^{-n(\eta y + |\eta|\delta) + n\Lambda(\eta)} = \tilde{\mu}_n(B_{y,\delta}) e^{-n\Lambda^*(y)} e^{-n|\eta|\delta} \end{aligned} \quad (2.2.53)$$

where the last equality follows by (2.2.45). Therefore, by combining (2.2.52) and (2.2.53)

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(B_{y,\delta}) \geq -\Lambda^*(y) - |\eta|\delta$$

Take now  $\delta_m \rightarrow 0$ ,  $\delta_m < \delta$ . Then, from the above,

$$\begin{aligned} \liminf_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(B_{y,\delta}) &\geq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(B_{y,\delta_m}) \\ &\geq -\Lambda^*(y) - |\eta|\delta_m \xrightarrow{m \rightarrow \infty} -\Lambda^*(y) \end{aligned}$$

and the proof of (2.2.51) is complete for  $y \in \mathcal{F}$ .

The inequality (2.2.51) trivially holds when  $x \notin \mathcal{D}_{\Lambda^\bullet}$ , while when  $x \in \mathcal{D}_{\Lambda^\bullet} \cap \overline{\mathcal{F}}$  there exist  $y_r \in \mathcal{F} \subseteq \mathcal{D}_{\Lambda^\bullet}$  such that  $y_r \rightarrow x$ . Since  $\Lambda^*(\cdot)$  is a convex rate function,  $\Lambda^*(y_r) \rightarrow \Lambda^*(x)$  and as  $B_{x,\delta}$  is an open set, eventually  $B_{y_r,\delta_r} \subseteq G$  for some  $\delta_r > 0$ . Finally, apply (2.2.51) for  $y_r \in \mathcal{F}$  to obtain

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(B_{x,\delta}) \geq \lim_{r \rightarrow \infty} \liminf_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(B_{y_r,\delta_r}) \geq - \lim_{r \rightarrow \infty} \Lambda^*(y_r) = -\Lambda^*(x).$$

(c+d). If  $\mathcal{D}_{\Lambda^\bullet} \subseteq \overline{\mathcal{F}}$  then for any  $G \subset \mathbb{R}^1$

$$- \inf_{x \in G \cap \overline{\mathcal{F}}} \Lambda^*(x) \geq - \inf_{x \in G \cap \mathcal{D}_{\Lambda^\bullet}} \Lambda^*(x) = - \inf_{x \in G} \Lambda^*(x)$$

The large deviations principle now follows by combining (2.2.46) and (2.2.47). The claim is a direct consequence of parts (b) and (c) in Lemma 2.2.2 below.  $\square$

**Remarks:**

- (a). The crucial step in the above proof of the upper bound is the exponential form of Markov's inequality combined with the independence assumption by which one can decompose the bounding exponent. For weakly dependent random variables a similar approach is to use the logarithmic limit of the left side of (2.2.50) instead of the logarithmic moment generating function for a single random variable. This is further explored in Section 2.3.
- (b). The crucial step in the above proof of the lower bound is the exponential change of measure applied when defining the associated measure  $\tilde{\mu}$ . This is particularly well suited to problems where, even if the random variables involved are not directly independent, some form of underlying independence exists (e.g., when a Girsanov type formula can be used, as in Chapter ??). Unless coupled with some convexity, this argument may in general fail in the dependent case. For more about it, c.f. Section ??.

The following lemma lifts some of the mystery behind the condition  $\mathcal{D}_{\Lambda^\bullet} \subseteq \overline{\mathcal{F}}$ .

**Lemma 2.2.2**

- (a).  $\mathcal{F} \subseteq \mathcal{D}_{\Lambda^\bullet}$  and both are intervals. Also  $0 \in \mathcal{D}_{\Lambda}$  which is an interval.
- (b). If  $0 \in \mathcal{D}_{\Lambda}^\circ$  and  $\Lambda$  is a steep function then  $\mathcal{D}_{\Lambda^\bullet} \subseteq \overline{\mathcal{F}}$ .
- (c). If  $0 \in \mathcal{D}_{\Lambda}^\circ$  then  $\Lambda^*$  is a good rate function.

**Remark:** Both the steepness of  $\Lambda$  and the condition  $0 \in \mathcal{D}_{\Lambda}^\circ$  are necessary for  $\mathcal{D}_{\Lambda^\bullet} \subseteq \overline{\mathcal{F}}$  (for details see exercise 2.2.3).

**Proof of lemma 2.2.1:**

(a). The convexity of  $\Lambda$  follows by Holder's inequality, as for any  $\theta \in [0, 1]$

$$\Lambda(\theta\lambda_1 + (1-\theta)\lambda_2) = \log E[(e^{\lambda_1 X_1})^\theta (e^{\lambda_2 X_1})^{(1-\theta)}] \leq \log \{E[e^{\lambda_1 X_1}]^\theta E[e^{\lambda_2 X_1}]^{(1-\theta)}\} = \theta\Lambda(\lambda_1) + (1-\theta)\Lambda(\lambda_2) .$$

Clearly  $\Lambda(0) = \log E[1] = 0$ , so  $\Lambda^*(x) \geq 0x - \Lambda(0) = 0$ . Let  $x_n \rightarrow x$ . Then for any  $\lambda \in \mathbb{R}$ ,

$$\liminf_{x_n \rightarrow x} \Lambda^*(x_n) \geq \liminf_{x_n \rightarrow x} [\lambda x_n - \Lambda(\lambda)] = \lambda x - \Lambda(\lambda)$$

and thus

$$\liminf_{x_n \rightarrow x} \Lambda^*(x_n) \geq \sup_{\lambda \in \mathbb{R}} [\lambda x - \Lambda(\lambda)] \triangleq \Lambda^*(x) ,$$

establishing the lower semicontinuity of  $\Lambda^*$ . Thus,  $\Lambda^*$  is a rate function.

The convexity of  $\Lambda^*$  follows by definition as

$$\begin{aligned} \theta\Lambda^*(x_1) + (1-\theta)\Lambda^*(x_2) &= \sup_{\lambda \in \mathbb{R}} \{\theta\lambda x_1 - \theta\Lambda(\lambda)\} + \sup_{\lambda \in \mathbb{R}} \{(1-\theta)\lambda x_2 - (1-\theta)\Lambda(\lambda)\} \\ &\geq \sup_{\lambda \in \mathbb{R}} \{(\theta x_1 + (1-\theta)x_2)\lambda - \Lambda(\lambda)\} = \Lambda^*(\theta x_1 + (1-\theta)x_2) . \end{aligned}$$

(b). By Jensen's inequality,

$$\Lambda(\lambda) = \log E[e^{\lambda X_1}] \geq E[\log e^{\lambda X_1}] = \lambda \bar{x} ,$$

for any  $\lambda \in \mathbb{R}$  and thus  $\Lambda^*(\bar{x}) = 0$  (this should have been expected in view of (2.2.39)).

Suppose now that  $x \in [\bar{x}, \infty)$ . Then, for any  $\lambda < 0$

$$\lambda x - \Lambda(\lambda) \leq \lambda \bar{x} - \Lambda(\lambda) \leq \Lambda^*(\bar{x}) = 0 ,$$

and (2.2.42) follows since  $\Lambda^*(x)$  is non-negative. Moreover,  $\Lambda^*(x)$  is nondecreasing for  $x > \bar{x}$  since for any  $\lambda \geq 0$  the function  $g(x) = \lambda x - \Lambda(\lambda)$  is nondecreasing. The proof of (2.2.43) for  $x \in (-\infty, \bar{x}]$  is similar.

(c). The identity (2.2.44) follows by interchanging the order of differentiation and integration. This is justified by the dominated convergence theorem since for  $\epsilon$  small enough

$$E[|X_1| e^{\eta X_1}] \leq \frac{1}{|\epsilon|} (M(\eta + \epsilon) + M(\eta - \epsilon)) < \infty .$$

Let  $\Lambda'(\eta) = y$  and consider the function  $g(\lambda) \triangleq \lambda y - \Lambda(\lambda)$ . As  $g(\cdot)$  is a concave function and  $g'(\eta) = 0$  clearly  $g(\eta) = \sup_{\lambda \in \mathbb{R}} g(\lambda)$  and (2.2.45) is established.  $\square$



**Proof of Lemma 2.2.2:**

(a). The sets  $\mathcal{D}_\Lambda$  and  $\mathcal{D}_{\Lambda^*}$  are convex since the functions  $\Lambda$  and  $\Lambda^*$  are convex. Further,  $\mathcal{D}_\Lambda$  and  $\mathcal{D}_{\Lambda^*}$  are intervals since any convex subset of  $\mathbb{R}^1$  is an interval. It was already noted that  $\Lambda(0) = 0$ , so  $0 \in \mathcal{D}_\Lambda$ . Moreover,  $\Lambda'(\cdot)$  is nondecreasing since  $\Lambda(\cdot)$  is convex. Now, since  $\mathcal{D}_\Lambda^o$  is an interval so is  $\mathcal{F}$ . Finally,  $\mathcal{F} \subseteq \mathcal{D}_{\Lambda^*}$  by (2.2.45).

(b). Since  $0 \in \mathcal{D}_\Lambda^o$  certainly  $\bar{x} = \Lambda'(0) \in \mathcal{F}$  and if  $\mathcal{D}_{\Lambda^*} = \{\bar{x}\}$  the proof is completed. Otherwise,  $\mathcal{D}_{\Lambda^*}^o$  is a non-empty interval. Consider now some  $x > \bar{x}$  and suppose that  $x \in \mathcal{D}_{\Lambda^*}^o$ . Recall that  $\Lambda^*(x) = \sup_{\lambda \in \mathcal{D}_\Lambda^+} g(\lambda)$  where the concave function  $g(\lambda) \triangleq \lambda x - \Lambda(\lambda)$  is continuous within the interval  $\mathcal{D}_\Lambda^+ \triangleq \mathcal{D}_\Lambda \cap [0, \infty)$  and differentiable in the non-empty interior of  $\mathcal{D}_\Lambda^+$ . Consequently,  $\Lambda^*(x) = \lim_{r \rightarrow \infty} g(\lambda_r)$  for some positive sequence  $\lambda_r \in \mathcal{D}_\Lambda^o$  such that  $g'(\lambda_r) \geq 0$ . Further, the sequence  $\{\lambda_r\}$  is bounded since  $\Lambda^*(x + \epsilon) < \infty$  for some  $\epsilon > 0$ , and as such has a limit point, say  $\lambda^*$ . Passing to the convergent subsequence,  $\lambda^* \notin \mathcal{D}_\Lambda^o$  implies by the steepness of  $\Lambda(\cdot)$  that  $\lim_{r \rightarrow \infty} g'(\lambda_r) = -\infty$ . This contradicts the above requirement that  $g'(\lambda_r) \geq 0$ . Therefore,  $\lambda^* \in \mathcal{D}_\Lambda^o$  implying  $x = \Lambda'(\lambda^*)$ , i.e.,  $x \in \mathcal{F}$ . A similar proof applies for  $x < \bar{x}$  so  $\mathcal{D}_{\Lambda^*}^o \subseteq \mathcal{F}$ . Since  $\mathcal{D}_{\Lambda^*}$  is an interval of non-empty interior it then follows that  $\mathcal{D}_{\Lambda^*} \subseteq \overline{\mathcal{F}}$ .

(c). If  $0 \in \mathcal{D}_\Lambda^o$  then there exist  $\lambda_+ > 0$  and  $\lambda_- < 0$  which are both in  $\mathcal{D}_\Lambda$ . Since for any  $\lambda \in \mathbb{R}^1$

$$\frac{\Lambda^*(x)}{|x|} \geq \lambda \text{sign}(x) - \frac{\Lambda(\lambda)}{|x|},$$

it follows that

$$\liminf_{|x| \rightarrow \infty} \frac{\Lambda^*(x)}{|x|} \geq \min\{\lambda_+, -\lambda_-\} > 0.$$

Thus, in particular  $\Lambda^*(x) \xrightarrow{|x| \rightarrow \infty} \infty$  and its level sets are bounded. Recall that closed and bounded subsets of  $\mathbb{R}^1$  are compact, so  $\Lambda^*$  is indeed a good rate function.  $\square$

**Exercises:**

**2.2.1** Suppose that  $\mathcal{D}_\Lambda = \mathbb{R}^1$ . Prove that  $\mu(\{x\}) = e^{-\Lambda^*(x)}$  when  $x \in \mathcal{D}_{\Lambda^*} \setminus \mathcal{F}$ .

**Hint:** Show first that there exists  $\{\lambda_\ell\}$  such that  $|\lambda_\ell| \rightarrow \infty$  and  $\int_{\mathbb{R}} e^{\lambda_\ell(z-x)} d\mu(z) = e^{-\Lambda^*(x)}$ .

**2.2.2** Explain why the above proof of (2.2.51) may not work when  $x \in \mathcal{D}_{\Lambda^*} \setminus \mathcal{F}$ .

**Hint:** Try distributions with density of the form  $Ce^{-\eta|x|}/(1 + |x|^p)$  for appropriate  $p$ .

2.2.3 (a). Prove that if  $\Lambda(\lambda) = \infty$  for all  $\lambda > 0$  and  $\bar{x} < \infty$ , then  $\Lambda^*(x) = 0$  for any  $x > \bar{x}$ , while  $\bar{\mathcal{F}} \subseteq (-\infty, \bar{x}]$ .

(b). Suppose now that  $\mathcal{D}_\Lambda = (-\infty, \lambda_0]$  where  $\lambda_0 > 0$  and  $\lim_{\lambda \uparrow \lambda_0} \Lambda'(\lambda) = x_0 < \infty$ . Prove that  $\Lambda^*(x) = \Lambda^*(x_0) + \lambda_0(x - x_0)$  for any  $x > x_0$  while  $\bar{\mathcal{F}} \subseteq (-\infty, x_0]$ .

2.2.4 Prove that  $\Lambda(\cdot)$  is lower semicontinuous.

Hint: Suppose that  $M(\lambda) = \infty$  for some  $\lambda > 0$  while  $\lambda - \delta \in \mathcal{D}_\Lambda^\circ$  for all  $\delta > 0$ . Let  $dG(x) \triangleq e^{\lambda x} d\mu(x)$  and observe that  $\lim_{\delta \rightarrow 0} G(\frac{1}{\delta}) = \infty$ . By integration by parts show that  $M(\lambda - \delta) \geq e^{-1} G(\frac{1}{\delta}) \rightarrow \infty$ . Conclude the proof using the convexity of  $\Lambda(\cdot)$ .

2.2.5 (a). Prove that  $\Lambda(\lambda)$  is  $C^\infty$  in  $\mathcal{D}_\Lambda^\circ$  and that  $\Lambda^*(x)$  is strictly convex and  $C^\infty$  in  $\mathcal{F}^\circ$ .

Hint: Show that  $x = \Lambda'(\eta) \in \mathcal{F}^\circ$  implies that  $\Lambda''(\eta) > 0$ .

(b). Construct an example where  $\mathcal{D}_\Lambda = \mathbb{R}^1$  while  $\Lambda^*(\cdot)$  is discontinuous.

Hint: Use a binary valued random variable.

2.2.6 Show that

(a) For  $X_1$  a Poisson( $\theta$ ) random variable  $\Lambda^*(x) = \theta - x + x \log(\frac{x}{\theta})$  for  $x \geq 0$  and  $\Lambda^*(x) = \infty$  otherwise.

(b) For  $X_1$  a Bernoulli( $p$ ) random variable  $\Lambda^*(x) = x \log(\frac{x}{p}) + (1-x) \log(\frac{1-x}{1-p})$  for  $x \in [0, 1]$  and  $\Lambda^*(x) = \infty$  otherwise.

(c) For  $X_1$  an Exponential ( $\theta$ ) random variable  $\Lambda^*(x) = \theta x - 1 - \log(\theta x)$  for  $x > 0$  and  $\Lambda^*(x) = \infty$  otherwise.

(d) For  $X_1$  a Normal  $(0, \sigma^2)$  random variable  $\Lambda^*(x) = \frac{x^2}{2\sigma^2}$ .

Verify that a large deviations principle holds in all these cases.

2.2.7 Let  $(x_-, x_+)$  be the largest interval around  $\bar{x}$  such that  $(x_-, x_+) \cap \bar{A} = \emptyset$ . Define  $I_A \triangleq \min\{\Lambda^*(x_-), \Lambda^*(x_+)\}$  and  $x^*$  any point at which this minimum is obtained.

(a). Prove that

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(A) = -I_A \quad (2.2.54)$$

whenever  $I_A = \infty$ .

(b). Prove that (2.2.54) also holds for  $I_A < \infty$  if  $x^* \in \overline{A^o \cap \mathcal{F}}$ .

**Hint:** Let  $x_r \in A^o \cap \mathcal{F}$  be a sequence which converges to  $x^*$  and apply (2.2.51) in order to verify that  $\liminf_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(A) \geq -\lim_{r \rightarrow \infty} \Lambda^*(x_r) = -\Lambda^*(x^*)$ .

**2.2.8** Assume that  $a \leq X_1 \leq b$

(a) Show that for any  $\lambda$

$$M(\lambda) \leq \frac{b - \bar{x}}{b - a} e^{\lambda a} + \frac{\bar{x} - a}{b - a} e^{\lambda b}.$$

This is the main ingredient of Hoeffding's inequality.

(b) Use (a) to show that for any  $a \leq x \leq b$

$$\Lambda^*(x) \geq H\left(\frac{x - a}{b - a} \middle| \frac{\bar{x} - a}{b - a}\right), \quad (2.2.55)$$

where  $H(p|p_0) \triangleq p \log(p/p_0) + (1 - p) \log((1 - p)/(1 - p_0))$ .

(c) Prove that the inequality (2.2.55) is sharp.

**2.2.9** Assume that  $\bar{x} = E[X_1]$  is finite,  $X_1 \leq b$  and  $\text{Var}(X_1) \leq \sigma^2$

(a) Show that for any  $\lambda > 0$

$$M(\lambda) \leq e^{\lambda \bar{x}} \left\{ \frac{(b - \bar{x})^2}{(b - \bar{x})^2 + \sigma^2} e^{-\lambda \sigma^2 / (b - \bar{x})} + \frac{\sigma^2}{(b - \bar{x})^2 + \sigma^2} e^{\lambda(b - \bar{x})} \right\}.$$

This is the main ingredient of Bennet's inequality.

(b) Use (a) to show that

$$\Lambda^*(x) \geq H(p_x | p_{\bar{x}}), \quad (2.2.56)$$

for any  $\bar{x} \leq x \leq b$  where  $p_x \triangleq \frac{(b - \bar{x})(x - \bar{x}) + \sigma^2}{(b - \bar{x})^2 + \sigma^2}$ .

(c) Prove that the inequality (2.2.56) is sharp.

## 2.3 A general large deviations principle in $\mathbb{R}^d$

Cramer's theorem 2.2.1 possesses a multivariate counterpart dealing with the large deviations of the empirical means of i.i.d real vectors in  $\mathbb{R}^d$ . Actually, the essential elements in the proof of this theorem extend to a more general class of dependent random vectors. This is explored here, where emphasis is put on the new points in which  $\mathbb{R}^d$  differs from  $\mathbb{R}^1$  with an eye towards possible infinite dimensional extensions.

Key applications are of-course Cramer's theorem in  $\mathbb{R}^d$  which is presented in Theorem 2.3.2 and its consequence – Sanov's Theorem for finite alphabets (see Corollary 2.3.1 and exercise 2.3.1). Some simple non i.i.d. applications are left as exercises 2.3.6, 2.3.8 and 2.3.9 while Section 2.4 is devoted to another class of key applications – the large deviations of Markov chains over finite alphabets.

The set-up considered here consists of a sequence of random vectors  $Z_n \in \mathbb{R}^d$  with laws  $\mu_n$  and logarithmic moment generating functions

$$\Lambda_n(\lambda) \triangleq \log E \left[ e^{\langle \lambda, Z_n \rangle} \right] \quad (2.3.57)$$

where  $\langle \lambda, x \rangle \triangleq \sum_{j=1}^d \lambda^j x^j$  is the usual scalar product in  $\mathbb{R}^d$  with  $x^j$  the  $j$ -th coordinate of the point  $x \in \mathbb{R}^d$  and  $|x| \triangleq \sqrt{\langle x, x \rangle}$  the usual Euclidean norm.

The existence of the limit of the properly scaled logarithmic moment generating functions indicates that  $\mu_n$  may satisfy a large deviations principle. Specifically, the following assumption prevails throughout this section.

**Assumption 2.3.1** *There exists a sequence of constants  $a_n \xrightarrow[n \rightarrow \infty]{} 0$  such that for each  $\lambda \in \mathbb{R}^d$ , the limit*

$$\Lambda(\lambda) \triangleq \lim_{n \rightarrow \infty} a_n \Lambda_n(a_n^{-1} \lambda) \quad (2.3.58)$$

*exists (possibly as a point in  $[-\infty, \infty]$ ). Further,  $0 \in \mathcal{D}_\Lambda^\circ$  where  $\mathcal{D}_\Lambda$  is the domain in which  $\Lambda(\cdot) < \infty$ .*

For example, when  $\mu_n$  is the law governing the empirical mean  $\hat{S}_n$  of the i.i.d. random vectors  $X_i \in \mathbb{R}^d$  then for any  $n \in \mathbb{Z}^+$

$$\frac{1}{n} \Lambda_n(n\lambda) = \Lambda(\lambda) \triangleq \log E[e^{\langle \lambda, X_1 \rangle}] \quad (2.3.59)$$

and the above assumption holds whenever  $0 \in \mathcal{D}_\Lambda^\circ$ .

**Definition 2.3.1** *The Fenchel-Legendre transform of  $\Lambda(\lambda)$  is*

$$\Lambda^*(x) \triangleq \sup_{\lambda \in \mathbb{R}^d} (\langle \lambda, x \rangle - \Lambda(\lambda)), \quad (2.3.60)$$

with  $\mathcal{D}_{\Lambda^*}$  denoting the domain in which  $\Lambda^*(\cdot) < \infty$ .

Cramer's theorem suggests that  $\Lambda^*$  is the natural candidate rate function for governing the large deviations principle associated with  $\mu_n$ . This is indeed proved in Theorem 2.3.1 under an additional condition. The properties of the functions  $\Lambda$  and  $\Lambda^*$  which are needed for that purpose are summarized in the following lemma whose proof is deferred to the end of this section.

**Lemma 2.3.1** *Assume 2.3.1.*

- (a).  $\Lambda(\lambda)$  is a convex function.  $\Lambda(\lambda) > -\infty$  everywhere and  $\Lambda^*(x)$  is a good, convex rate function.
- (b).  $\mathcal{D}_\Lambda$  and  $\mathcal{D}_{\Lambda^*}$  are convex sets.
- (c). Suppose that  $y = \nabla \Lambda(\eta)$  for some  $\eta \in \mathcal{D}_\Lambda^\circ$ . Then

$$\Lambda^*(y) = \langle \eta, y \rangle - \Lambda(\eta). \quad (2.3.61)$$

- (d). Let  $y, \eta$  be as in (c) above. Then, for any  $x \neq y$

$$\Lambda_\eta^*(x) > \Lambda_\eta^*(y) = 0, \quad (2.3.62)$$

where  $\Lambda_\eta^*(\cdot)$  is the Fenchel-Legendre transform of

$$\Lambda_\eta(\theta) \triangleq \Lambda(\theta + \eta) - \Lambda(\eta). \quad (2.3.63)$$

**Remark:** These convex analysis considerations are addressed again in Section ?? in a more abstract setup.

The general large deviations principle in  $\mathbb{R}^d$  is now summarized as follows.

**Theorem 2.3.1 (Gartner)** *Assume 2.3.1. Then*

- (a). For any closed set  $F$

$$\limsup_{n \rightarrow \infty} a_n \log \mu_n(F) \leq - \inf_{x \in F} \Lambda^*(x). \quad (2.3.64)$$

(b). For any open set  $G$

$$\liminf_{n \rightarrow \infty} a_n \log \mu_n(G) \geq - \inf_{x \in G \cap \overline{\mathcal{F}}} \Lambda^*(x), \quad (2.3.65)$$

where  $\mathcal{F} \triangleq \{x : x = \nabla \Lambda(\lambda) \text{ for some } \lambda \in \mathcal{D}_\Lambda^\circ\}$ .

(c). If  $\mathcal{D}_{\Lambda^\bullet} \subseteq \overline{\mathcal{F}}$  then the family of probability measures  $\mu_n$  satisfies a large deviations principle controlled by the good rate function  $\Lambda^*(x)$ .

The following lemma, whose proof is also deferred to the end of this section, makes the above theorem applicable by stating explicit conditions on  $\Lambda$  under which  $\mathcal{D}_{\Lambda^\bullet} \subseteq \overline{\mathcal{F}}$ .

**Lemma 2.3.2 (Ellis)** Suppose that  $\Lambda(\lambda)$  which satisfies 2.3.1 is a lower semicontinuous function which is differentiable in  $\mathcal{D}_\Lambda^\circ$  and moreover it is a steep function, namely for any convergent sequence  $\lambda_n \in \mathcal{D}_\Lambda^\circ$  whose limit does not belong to  $\mathcal{D}_\Lambda^\circ$ ,  $\lim_{n \rightarrow \infty} |\nabla \Lambda(\lambda_n)| = \infty$ . Then,  $\mathcal{D}_{\Lambda^\bullet} \subseteq \overline{\mathcal{F}}$ .

The proof of Theorem 2.3.1 is given in the sequel, preceded by the statement and derivation of a key application – Cramer’s theorem about the empirical means of i.i.d. random vectors in  $\mathbb{R}^d$ .

**Theorem 2.3.2** Let  $\mu_n$  be the laws governing the empirical means  $\hat{S}_n \triangleq \frac{1}{n} \sum_{i=1}^n X_i$  where  $X_i \in \mathbb{R}^d$  are i.i.d. distributed according to the law  $\mu$ . Suppose that the logarithmic moment generating function

$$\Lambda(\lambda) \triangleq \log E_\mu [e^{\langle \lambda, X_1 \rangle}] , \quad (2.3.66)$$

is a steep lower semicontinuous function which is finite in some open ball centered at the origin (in particular these conditions hold when  $\mathcal{D}_\Lambda = \mathbb{R}^d$ ). Then, for any measurable set  $\Gamma \subset \mathbb{R}^d$

$$- \inf_{x \in \Gamma^\circ} \Lambda^*(x) \leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(\Gamma) \leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log \mu_n(\Gamma) \leq - \inf_{x \in \overline{\Gamma}} \Lambda^*(x) , \quad (2.3.67)$$

with  $\Lambda^*$  being the Fenchel-Legendre transform of  $\Lambda$ .

**Proof:** Recall that in this case the basic assumption 2.3.1 holds (see (2.3.59)) and indeed  $\Lambda$  of (2.3.58) is given by (2.3.66). In the statement of the theorem it is further assumed that  $\Lambda$  is a steep function. It follows from the definition (2.3.66) that  $\Lambda$  is differentiable in  $\mathcal{D}_\Lambda^\circ$  (by an argument which is similar to the proof of (2.2.44) in Section 2.2). Thus, Lemma 2.3.2 applies and (2.3.67) follows by Theorem 2.3.1.  $\square$

**Remark** The assumption of lower semi-continuity may be relaxed. For details, c.f. exercise 2.3.3.

The essential ingredients for the proof of Theorem 2.3.1 are those presented in Section 2.2 in the course of proving Theorem 2.2.1, namely – the exponential form of Markov's inequality is applied for deriving the upper bound and an exponential change of measure is used for deriving the lower bound. However, here one encounters two new obstacles which slightly complicate both parts of the proof.

**Proof of Theorem 2.3.1:**

(a). In  $\mathbb{R}^d$  the monotonicity of  $\Lambda^*$  stated in Lemma 2.2.1 part (b) is somewhat lost. Thus, the strategy of containing  $\bar{A}$  by two half-spaces  $A_-$  and  $A_+$  is not as useful as it is in  $\mathbb{R}^1$ . Instead, here one uses Markov's inequality to obtain tight upper bounds for all small closed balls. Then, compact sets are covered by an appropriate finite collection of small enough balls and the upper bound follows for compact sets by the union of events bound.

As mentioned in Section 1.1, proving (2.3.64) is equivalent to proving that for any  $\delta > 0$  and any closed set  $F \subset \mathbb{R}^d$

$$\limsup_{n \rightarrow \infty} a_n \log \mu_n(F) \leq \delta - \inf_{x \in F} I^\delta(x) \quad (2.3.68)$$

where

$$I^\delta(x) \triangleq \begin{cases} \Lambda^*(x) - \delta & x \in \mathcal{D}_\Lambda \\ \frac{1}{\delta} & x \notin \mathcal{D}_\Lambda \end{cases}$$

Fix now  $\delta > 0$  and an arbitrary compact set  $\Gamma$ . For any  $q \in \Gamma$  choose  $\lambda_q \in \mathcal{D}_\Lambda$  for which

$$\langle \lambda_q, q \rangle - \Lambda(\lambda_q) \geq I^\delta(q) \quad (2.3.69)$$

This is always possible by the definitions of  $\Lambda^*$  and  $I^\delta$ . Choose now  $\rho_q > 0$  such that  $\rho_q |\lambda_q| \leq \delta$  and let  $B_{q, \rho_q}$  be the open ball with center at the point  $q$  and radius  $\rho_q$  with  $\bar{B}_{q, \rho_q}$  the corresponding closed ball.

Then, for any  $n$  and any  $q \in \Gamma$

$$\mu_n(B_{q, \rho_q}) \leq \exp \left( - \inf_{x \in \bar{B}_{q, \rho_q}} \{a_n^{-1} \langle \lambda_q, x \rangle\} \right) E \left[ \exp(a_n^{-1} \langle \lambda_q, Z_n \rangle) \right].$$

Thus,

$$a_n \log \mu_n(B_{q, \rho_q}) \leq - \inf_{x \in \bar{B}_{q, \rho_q}} \langle \lambda_q, x \rangle + a_n \Lambda_n(a_n^{-1} \lambda_q) \leq \delta - \langle \lambda_q, q \rangle + a_n \Lambda_n(a_n^{-1} \lambda_q). \quad (2.3.70)$$

As  $\Gamma$  is a compact set, one can extract from the open cover  $\bigcup_{q \in \Gamma} B_{q, \rho_q}$  of  $\Gamma$  a finite cover which consists of  $N < \infty$  (depending only on  $\Gamma$  and  $\delta$ ) such balls with centers  $q_1, \dots, q_N$  in  $\Gamma$ . By the union of events bound and (2.3.70),

$$a_n \log \mu_n(\Gamma) \leq a_n \log N + \delta - \min_{i=1, \dots, N} \{ \langle \lambda_{q_i}, q_i \rangle - a_n \Lambda_n(a_n^{-1} \lambda_{q_i}) \}.$$

Since  $a_n \rightarrow 0$  as  $n \rightarrow \infty$  while  $a_n \Lambda_n(a_n^{-1} \lambda_{q_i}) \rightarrow \Lambda(\lambda_{q_i})$  (uniformly over  $i = 1, \dots, N$ ), one obtains

$$\limsup_{n \rightarrow \infty} a_n \log \mu_n(\Gamma) \leq \delta - \min_{i=1, \dots, N} \{ \langle \lambda_{q_i}, q_i \rangle - \Lambda(\lambda_{q_i}) \} \leq \delta - \min_{i=1, \dots, N} I^\delta(q_i),$$

where the last inequality follows from (2.3.69). As  $q_i \in \Gamma$  the upper bound (2.3.68) is thus established for all compact sets.

This upper bound is extended to all closed sets in  $\mathbb{R}^d$  by showing that  $\mu_n$  is an *exponentially tight* family of probability measures and using Lemma 1.1.1. Specifically, it is shown in the sequel that for any  $\alpha < \infty$  there exists  $\rho_\alpha$  large enough such that for the compact set  $K_\alpha \triangleq \overline{B}_{0, \rho_\alpha}$

$$\limsup_{n \rightarrow \infty} a_n \log \mu_n(K_\alpha^c) < -\alpha \quad (2.3.71)$$

For that purpose, observe first that  $\overline{B}_{0, \rho d}^c \subset \bigcup_{j=1}^d \{x : |x^j| \geq \rho\}$ . Therefore, by the union of events bound

$$\begin{aligned} \mu_n(\overline{B}_{0, \rho d}^c) &\leq \mu_n(\bigcup_{j=1}^d \{x : |x^j| \geq \rho\}) \\ &\leq \sum_{j=1}^d \{ \mu_n^j([\rho, \infty)) + \mu_n^j((-\infty, -\rho]) \} \end{aligned} \quad (2.3.72)$$

where  $\mu_n^j$  are the laws of  $Z_n^j$ , the coordinates of the random vector  $Z_n$ . As  $\mathcal{D}_\Lambda$  contains an open ball around the origin there exist  $\theta_j^+ > 0$  and  $\theta_j^- < 0$  such that  $\Lambda(\theta_j^+ \mathbf{u}_j) < \infty$  and  $\Lambda(\theta_j^- \mathbf{u}_j) < \infty$  where  $\mathbf{u}_j$  denotes the  $j$ -th unit vector in  $\mathbb{R}^d$  for  $j = 1, \dots, d$ . By the exponential form of Markov's inequality, for  $j = 1, \dots, d$

$$\limsup_{n \rightarrow \infty} a_n \log \mu_n^j([\rho, \infty)) \leq -\theta_j^+ \rho + \limsup_{n \rightarrow \infty} a_n \Lambda_n(a_n^{-1} \theta_j^+ \mathbf{u}_j) = -\theta_j^+ \rho + \Lambda(\theta_j^+ \mathbf{u}_j). \quad (2.3.73)$$

Similarly,

$$\limsup_{n \rightarrow \infty} a_n \log \mu_n^j((-\infty, -\rho]) \leq \theta_j^- \rho + \Lambda(\theta_j^- \mathbf{u}_j). \quad (2.3.74)$$



Now, the inequality (2.3.71) results by combining (2.3.72), (2.3.73) and (2.3.74) and considering  $\rho \rightarrow \infty$ .

(b). Focusing now on establishing the lower bound (2.3.65) for any open set, it suffices to prove

$$\lim_{\delta \rightarrow 0} \liminf_{n \rightarrow \infty} a_n \log \mu_n(B_{y,\delta}) \geq -\Lambda^*(y), \quad (2.3.75)$$

for any  $y \in \mathcal{F}$ . Indeed, (2.3.75) then holds for any  $y \in \overline{\mathcal{F}}$  and (2.3.65) follows by the same argument encountered in that leads from (2.2.47) to (2.2.51) in the proof of Theorem 2.2.1.

Fix now  $y = \nabla \Lambda(\eta) \in \mathcal{F}$  with  $\eta \in \mathcal{D}_\Lambda^0$ . Then, for all  $n$  large enough,  $\Lambda_n(a_n^{-1}\eta) < \infty$  and the “associated” probability measures

$$d\tilde{\mu}_n(z) = \exp[a_n^{-1} \langle \eta, z \rangle - \Lambda_n(a_n^{-1}\eta)] d\mu_n(z) \quad (2.3.76)$$

are well defined. Clearly,

$$\begin{aligned} a_n \log \mu_n(B_{y,\delta}) &= a_n \Lambda_n(a_n^{-1}\eta) - \langle \eta, y \rangle + a_n \log \int_{z \in B_{y,\delta}} \exp(a_n^{-1} \langle \eta, y - z \rangle) d\tilde{\mu}_n(z) \\ &\geq a_n \Lambda_n(a_n^{-1}\eta) - \langle \eta, y \rangle - |\eta|\delta + a_n \log \tilde{\mu}_n(B_{y,\delta}) \end{aligned} \quad (2.3.77)$$

Therefore,

$$\begin{aligned} \lim_{\delta \rightarrow 0} \liminf_{n \rightarrow \infty} a_n \log \mu_n(B_{y,\delta}) &\geq \Lambda(\eta) - \langle \eta, y \rangle + \lim_{\delta \rightarrow 0} \liminf_{n \rightarrow \infty} a_n \log \tilde{\mu}_n(B_{y,\delta}) \\ &= -\Lambda^*(y) + \lim_{\delta \rightarrow 0} \liminf_{n \rightarrow \infty} a_n \log \tilde{\mu}_n(B_{y,\delta}) \end{aligned} \quad (2.3.78)$$

where the above equality follows by (2.3.61).

Here, a new obstacle stems from the removal of the independence assumption. Indeed,  $E_{\tilde{\mu}_n}(Z_n) \rightarrow y$  as  $n \rightarrow \infty$ , but now one still has to establish the appropriate analog of the weak law of large numbers. This is handled in the sequel by applying the *large deviations upper bound* for the “associated” family of measures  $\tilde{\mu}_n$ . Indeed, the proof of (2.3.75) is completed by showing that for any  $\delta > 0$

$$\limsup_{n \rightarrow \infty} a_n \log \tilde{\mu}_n(B_{y,\delta}^c) < 0. \quad (2.3.79)$$

For that purpose let  $\Lambda_{n,\eta}(\cdot)$  denotes the logarithmic moment generating function corresponding to the law  $\tilde{\mu}_n$ . Then, for any  $\theta \in \mathbb{R}^d$

$$a_n \Lambda_{n,\eta}(a_n^{-1}\theta) \triangleq a_n \log \left[ \int_{\mathbb{R}^d} e^{a_n^{-1} \langle \theta, z \rangle} d\tilde{\mu}_n(z) \right] = a_n \Lambda_n(a_n^{-1}(\theta + \eta)) - a_n \Lambda_n(a_n^{-1}\eta) \rightarrow \Lambda_\eta(\theta)$$

as  $a_n \Lambda_n(a_n^{-1} \eta) < \infty$  for all  $n$  large (recall that  $\eta \in \mathcal{D}_\Lambda$ ). Moreover,  $\Lambda_\eta(0) = 0$  and  $\Lambda_\eta(\theta) < \infty$  for all  $|\theta|$  small enough since  $\eta \in \mathcal{D}_\Lambda^\circ$  (c.f. (2.3.63)). Thus, a large deviations *upper bound* of the form of (2.3.64) holds for the sequence of measures  $\tilde{\mu}_n$ . In particular, for the closed set  $B_{y,\delta}^c$  it yields

$$\limsup_{n \rightarrow \infty} a_n \log \tilde{\mu}_n(B_{y,\delta}^c) \leq - \inf_{x \in B_{y,\delta}^c} \Lambda_\eta^*(x). \quad (2.3.80)$$

Moreover,  $\Lambda_\eta^*(x) > 0$  for any  $x \neq y$  in view of (2.3.62) and paralleling the proof of part (a) of Lemma 2.3.1 one easily shows that  $\Lambda_\eta^*$  is a good rate function. Therefore,  $\inf_{x \in B_{y,\delta}^c} \Lambda_\eta^*(x) > 0$  for all  $\delta > 0$  and (2.3.80) implies (2.3.79), concluding the proof of the lower bound (2.3.65).

(c). The large deviations principle follows by combining (2.3.64) and (2.3.65) as  $\mathcal{D}_{\Lambda^\bullet} \subseteq \overline{\mathcal{F}}$  implies that for any  $G$

$$- \inf_{x \in G \cap \overline{\mathcal{F}}} \Lambda^*(x) \geq - \inf_{x \in G \cap \mathcal{D}_{\Lambda^\bullet}} \Lambda^*(x) = - \inf_{x \in G} \Lambda^*(x)$$

□

#### Remarks:

- (a). The proof above actually extends beyond  $\mathbb{R}^d$  and in principle is applicable for any metric space (as shown in Chapter ??). However, two points of caution are that the exponential tightness has to be proved on a case by case basis and that in infinite dimensional spaces  $\Lambda$  is rarely differentiable and therefore the analog of Lemma 2.3.2 typically fails to hold.
- (b). In the current set-up the condition  $0 \in \mathcal{D}_\Lambda^\circ$  which is needed for the above proof does not follow as a consequence of  $\mathcal{D}_{\Lambda^\bullet} \subseteq \overline{\mathcal{F}}$ . For this reason it is incorporated in Assumption 2.3.1. Note however that the condition  $0 \in \mathcal{D}_\Lambda^\circ$  is not required at all for proving the upper bound (2.3.64) for compact sets.

As mentioned in Section 2.1, Sanov's theorem for finite alphabets, Theorem 2.1.1, may be deduced as a consequence of Cramer's theorem 2.3.2. Indeed, note that the empirical mean of the random vectors  $Y_i \triangleq [1_{X_i=a_1}, 1_{X_i=a_2}, \dots, 1_{X_i=a_{|\Sigma|}}]$  equals  $L_n^X$ , the empirical measure of the i.i.d. random variables  $X_1, \dots, X_n$  over the finite alphabet  $\Sigma$ . Moreover, as  $Y_i$  are bounded, here  $\mathcal{D}_\Lambda = \mathbb{R}^{|\Sigma|}$  and one obtains the following corollary of Theorem 2.3.2.

**Corollary 2.3.1** *For any set  $\Gamma$  of probability vectors in  $\mathbb{R}^{|\Sigma|}$*

$$- \inf_{\nu \in \Gamma^\circ} \Lambda^*(\nu) \leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \text{Prob}_\mu(L_n^X \in \Gamma) \leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log \text{Prob}_\mu(L_n^X \in \Gamma) \leq - \inf_{\nu \in \Gamma} \Lambda^*(\nu) \quad (2.3.81)$$

where  $\Lambda^*$  is the Fenchel-Legendre transform of the logarithmic moment generating function

$$\Lambda(\lambda) = \log E_\mu(e^{<\lambda, Y_1>}) = \log \sum_{i=1}^{|\Sigma|} e^{\lambda_i} \mu(a_i) \quad (2.3.82)$$

and  $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_{|\Sigma|}) \in \mathbb{R}^{|\Sigma|}$ .

**Remark:** Comparing the above corollary to Theorem 2.1.1 it is tempting to conjecture that  $\Lambda^*(\cdot) = H(\cdot|\mu)$ . Indeed, this is proved in exercise 2.3.1. Actually, as shown in Section ?? the rate function controlling a large deviations principle in  $\mathbb{R}^d$  is always unique, thus this result is not surprising.

**Proof of Lemma 2.3.1:**

(a). Since  $\Lambda_n$  are convex functions (see the proof of part (a) of Lemma 2.2.1) so are  $a_n \Lambda_n(a_n^{-1} \cdot)$  and their limit  $\Lambda(\cdot)$  is convex as well. Moreover,  $\Lambda_n(0) = 0$  and therefore  $\Lambda(0) = 0$  implying that  $\Lambda^*$  is non-negative. Both the convexity and the lower semicontinuity of  $\Lambda^*$  follow from its definition via (2.3.60) (see proof of part (a) of Lemma 2.2.1).

Now if  $\Lambda(\lambda) = -\infty$  for some  $\lambda \in \mathbb{R}^d$  then by convexity  $\Lambda(\alpha\lambda) = -\infty$  for all  $\alpha \in (0, 1]$ . Moreover, since  $\Lambda(0) = 0$  it follows by convexity that  $\Lambda(-\alpha\lambda) = \infty$  for all  $\alpha \in (0, 1]$  contradicting the assumption that  $0 \in \mathcal{D}_\Lambda^o$ .

Since  $0 \in \mathcal{D}_\Lambda^o$ , it follows that  $\overline{B}_{0,\delta} \subset \mathcal{D}_\Lambda^o$  for some  $\delta > 0$  and  $C = \sup_{\lambda \in \overline{B}_{0,\delta}} \Lambda(\lambda) < \infty$  since the convex function  $\Lambda$  is continuous in  $\mathcal{D}_\Lambda^o$ . Therefore,

$$\Lambda^*(x) \geq \sup_{\lambda \in \overline{B}_{0,\delta}} \{<\lambda, x> - \Lambda(\lambda)\} \geq \sup_{\lambda \in \overline{B}_{0,\delta}} <\lambda, x> - \sup_{\lambda \in \overline{B}_{0,\delta}} \Lambda(\lambda) = \delta|x| - C \quad (2.3.83)$$

Thus, the level sets  $\{x : \Lambda^*(x) \leq \alpha\}$  are clearly bounded within a closed ball around the origin (of radius  $(\alpha + C)/\delta$ ) and  $\Lambda^*$  is necessarily a good rate function.

(b). The convexity of  $\mathcal{D}_\Lambda$  and  $\mathcal{D}_{\Lambda^*}$  is merely a consequence of the convexity of  $\Lambda$  and  $\Lambda^*$  respectively.

(c). Clearly  $\Lambda^*(y) \geq [<\eta, y> - \Lambda(\eta)]$ . Assume that this inequality is strict, i.e., for some  $\lambda \in \mathcal{D}_\Lambda$

$$[<\lambda, y> - \Lambda(\lambda)] > [<\eta, y> - \Lambda(\eta)]. \quad (2.3.84)$$

Since  $\mathcal{D}_\Lambda$  is a convex set,

$$g(\alpha) \triangleq \alpha <\lambda - \eta, y> - \Lambda(\eta + \alpha(\lambda - \eta)) + [<\eta, y> - \Lambda(\eta)] \quad \alpha \in [0, 1] \quad (2.3.85)$$

is a finite valued concave function. Thus, by concavity

$$g(1) - g(0) \leq \liminf_{\alpha \downarrow 0} \frac{g(\alpha) - g(0)}{\alpha} = \langle \lambda - \eta, y - \nabla \Lambda(\eta) \rangle = 0, \quad (2.3.86)$$

where the last equality follows from the assumption  $y = \nabla \Lambda(\eta)$ . However, the inequality above contradicts (2.3.84) and therefore (2.3.61) is established.

(d). The function  $\Lambda_\eta^*$  is non-negative (being the Fenchel-Legendre transform of  $\Lambda_\eta$  where  $\Lambda_\eta(0) = 0$ ) and moreover

$$\Lambda_\eta^*(x) = \sup_{\theta \in \mathbb{R}^d} \{ \langle \theta, x \rangle - \Lambda(\theta + \eta) \} + \Lambda(\eta) = \Lambda^*(x) - \langle \eta, x \rangle + \Lambda(\eta). \quad (2.3.87)$$

Thus,  $\Lambda_\eta^*(y) = 0$  by (2.3.61) and  $\Lambda_\eta^*(x) \geq \Lambda_\eta^*(y)$  for any  $x \neq y$ . If  $\Lambda_\eta^*(x) = 0$  for some  $x$  then for any  $\theta \in \mathbb{R}^d$

$$\langle \theta, x \rangle \leq \Lambda_\eta(\theta). \quad (2.3.88)$$

In particular also

$$\langle \theta, x \rangle \leq \lim_{n \rightarrow \infty} \frac{\Lambda_\eta^*(\theta/n)}{(1/n)} = \langle \theta, \nabla \Lambda_\eta(0) \rangle, \quad (2.3.89)$$

where  $\nabla \Lambda_\eta(0) = \nabla \Lambda(\eta) = y$ . Since the inequality (2.3.89) holds for all  $\theta \in \mathbb{R}^d$ , necessarily  $x = \nabla \Lambda_\eta(0) = y$ . Thus,  $\Lambda_\eta^*(x) > 0$  for any  $x \neq y$  and (2.3.62) is established.  $\square$

#### Proof of Lemma 2.3.2:

Since  $\mathcal{D}_\Lambda^\circ$  is not empty, by (2.3.61) so is  $\mathcal{D}_{\Lambda^\bullet}$ . Further, if  $\mathcal{D}_{\Lambda^\bullet} = \{x\}$  then necessarily  $x = \nabla \Lambda(0) \in \mathcal{F}$  and the proof is then complete. Therefore, we may assume from here on that the convex set  $\mathcal{D}_{\Lambda^\bullet}$  contains at least one line so its **relative interior**

$$\text{ri } \mathcal{D}_{\Lambda^\bullet} \triangleq \{x \in \mathcal{D}_{\Lambda^\bullet} : y \in \mathcal{D}_{\Lambda^\bullet} \Rightarrow x - \epsilon(y - x) \in \mathcal{D}_{\Lambda^\bullet} \text{ for some } \epsilon > 0\}$$

is non-empty and moreover

$$\mathcal{D}_{\Lambda^\bullet} = \overline{\text{ri } \mathcal{D}_{\Lambda^\bullet}}. \quad (2.3.90)$$

The subdifferential  $\partial \Lambda(\lambda)$  of the convex function  $\Lambda(\cdot)$  at a point  $\lambda \in \mathbb{R}^d$ , is defined as

$$\partial \Lambda(\lambda) \triangleq \{x : \Lambda(\theta) \geq \Lambda(\lambda) + \langle x, \theta - \lambda \rangle \quad \forall \theta \in \mathbb{R}^d\} = \{x : \langle \lambda, x \rangle - \Lambda^*(x) \geq \Lambda(\lambda)\}, \quad (2.3.91)$$

where the second equality is a direct consequence of the definition of the Fenchel-Legendre transform. It is easy to check that  $\partial \Lambda(\lambda) = \emptyset$  for  $\lambda \notin \mathcal{D}_\Lambda$  (recall that  $\Lambda > -\infty$  everywhere) while

otherwise  $\partial\Lambda(\eta) = \{x : \Lambda_\eta^*(x) = 0\}$  (where  $\Lambda_\eta^*$  is defined in (2.3.63)). Since  $\Lambda$  is differentiable throughout  $\mathcal{D}_\Lambda^\circ$  (by assumption), it follows by (2.3.61) and (2.3.62) that  $\partial\Lambda(\eta) = \{\nabla\Lambda(\eta)\}$  for any  $\eta \in \mathcal{D}_\Lambda^\circ$ .

The proof of the lemma is now divided into the following two steps:

- (a). Since  $\Lambda$  is lower semicontinuous, for any  $x \in \text{ri } \mathcal{D}_{\Lambda^\bullet}$  there exists  $\lambda \in \mathbb{R}^d$  such that  $x \in \partial\Lambda(\lambda)$ .
- (b). Since  $\Lambda$  is a steep function,  $\partial\Lambda(\lambda) = \emptyset$  for any  $\lambda \in \mathcal{D}_\Lambda \setminus \mathcal{D}_\Lambda^\circ$ .

When combined, these two claims result with

$$\text{ri } \mathcal{D}_{\Lambda^\bullet} \subseteq \bigcup_{\lambda \in \mathbb{R}^d} \partial\Lambda(\lambda) = \bigcup_{\lambda \in \mathcal{D}_\Lambda^\circ} \{\nabla\Lambda(\lambda)\} = \mathcal{F}, \quad (2.3.92)$$

which, together with (2.3.90), amount to the proof of the lemma.

- (a). Fix a point  $x \in \text{ri } \mathcal{D}_{\Lambda^\bullet}$  and define the function

$$g(y) \triangleq \inf_{\delta > 0} \frac{\Lambda^*(x + \delta y) - \Lambda^*(x)}{\delta} = \lim_{\delta \downarrow 0} \frac{\Lambda^*(x + \delta y) - \Lambda^*(x)}{\delta}, \quad (2.3.93)$$

where the convexity of  $\Lambda^*$  results with a monotonicity in  $\delta$  which in turn implies the above equality (and that the above limit exists). For the same reason  $g(y)$  is a convex function and the set  $\mathcal{E}_g \triangleq \{(y, \xi) : \xi \geq g(y)\} \subseteq \mathbb{R}^d \times \mathbb{R}$  is a convex set. Further,  $g(\alpha y) = \alpha g(y)$  for all  $\alpha \geq 0$  and in particular  $g(0) = 0$ . Observe that  $g(y) = \infty$  when  $x + \delta y \notin \mathcal{D}_{\Lambda^\bullet}$  for all  $\delta > 0$ . Consider therefore those directions  $y$  such that  $g(y) < \infty$ . Since  $x \in \text{ri } \mathcal{D}_{\Lambda^\bullet}$ , it then follows that for some  $\epsilon > 0$  the whole line segment  $x + \beta y$  for  $|\beta| \leq \epsilon$  is in  $\mathcal{D}_{\Lambda^\bullet}$ . Let  $\hat{y} \triangleq \epsilon y$  so by the convexity of  $\Lambda^*$

$$\Lambda^*(x) < (1 - \delta)\Lambda^*(x + \delta\hat{y}) + \delta\Lambda^*(x - (1 - \delta)\hat{y}) < \infty \quad (2.3.94)$$

for all  $\delta \in [0, 1]$ . This implies

$$g(\hat{y}) \geq \lim_{\delta \downarrow 0} [\Lambda^*(x + \delta\hat{y}) - \Lambda^*(x - (1 - \delta)\hat{y})] \geq -\Lambda^*(x - \hat{y}) = -\Lambda^*(x - \epsilon y) > -\infty \quad (2.3.95)$$

where the last two inequalities follow by the non-negativity and upper semicontinuity of  $\Lambda^*$ , and the fact that  $x - \epsilon y \in \mathcal{D}_{\Lambda^\bullet}$ . Since  $g(y) = \frac{g(\hat{y})}{\epsilon}$  it follows from the above that  $g(y) > -\infty$  for all  $y$ . Since  $g(y) = |y|g(y/|y|)$  and  $\inf_{\{z: |z|=1\}} g(z) > -\infty$  it now follows that  $\liminf_{y \rightarrow 0} g(y) \geq 0$ , implying that  $(0, -1) \notin \overline{\mathcal{E}}_g$ . The set  $\overline{\mathcal{E}}_g$  is closed, convex and non-empty (for example  $(0, 1) \in \overline{\mathcal{E}}_g$ ). Thus, there

exists a hyperplane in  $\mathbb{R}^d \times \mathbb{R}$  which strictly separates the point  $(0, -1)$  and the set  $\mathcal{E}_g$  (this is a particular instance of the Hahn-Banach theorem quoted in Appendix ??). Specifically, there exist  $\lambda \in \mathbb{R}^d$  and  $\rho \in \mathbb{R}$  such that

$$\langle \lambda, 0 \rangle + \rho = \rho > \langle \lambda, y \rangle - \xi \rho \quad \forall (y, \xi) \in \mathcal{E}_g. \quad (2.3.96)$$

Considering  $y = 0$  it is clear that  $\rho > 0$  and then by specializing (2.3.96) to  $\xi = g(y)$  one obtains that  $g^*(\frac{\lambda}{\rho}) \leq 1$  where  $g^*$  is the Fenchel-Legendre transform of  $g(y)$ . Observe now that  $g^*$  assumes only the values 0 or  $\infty$  as

$$g^*(\lambda) \triangleq \sup_y (\langle \lambda, y \rangle - g(y)) = \sup_{\alpha \geq 0} \alpha \left[ \sup_{\{z: |z|=1\}} (\langle \lambda, z \rangle - g(z)) \right] = \begin{cases} 0 & \lambda \in C \\ \infty & \lambda \notin C \end{cases} \quad (2.3.97)$$

where  $C = \{\lambda : |z| = 1 \Rightarrow g(z) \geq \langle \lambda, z \rangle\}$ . Thus, the set  $C$  must be non-empty or equivalently there exists  $\lambda_0 \in \mathbb{R}^d$  such that  $g(y) \geq \langle \lambda_0, y \rangle$  for all  $y \in \mathbb{R}^d$ .

Considering now (2.3.93) one obtains

$$\Lambda^*(z) - \Lambda^*(x) \geq \langle \lambda_0, z - x \rangle$$

for all  $z \in \mathbb{R}^d$ . Therefore, also

$$\langle \lambda_0, x \rangle - \Lambda^*(x) = \sup_{z \in \mathbb{R}^d} (\langle \lambda_0, z \rangle - \Lambda^*(z)) \geq \Lambda(\lambda_0), \quad (2.3.98)$$

where the above inequality is a property of the Legendre transform of any convex, lower semicontinuous function  $f$  such that  $f(\cdot) > -\infty$  everywhere and  $f(\cdot) \not\equiv \infty$  (consider Section ?? for a proof). Indeed, it is assumed in the statement of this lemma that  $\Lambda(\cdot)$  is lower semicontinuous and the other conditions mentioned above are satisfied in view of part (a) of Lemma 2.3.1. The inequality (2.3.98) amounts to  $x \in \partial\Lambda(\lambda_0)$  for some  $\lambda_0 \in \mathcal{D}_\Lambda$ .

(b). Suppose there exists a point  $\eta \in \mathcal{D}_\Lambda \setminus \mathcal{D}_\Lambda^o$  such that  $\partial\Lambda(\eta)$  is non-empty. Then,  $\partial\Lambda(\eta) = \{x : \Lambda_\eta^*(x) = 0\}$  is a non-empty closed convex set (since  $\Lambda_\eta^*$  is a convex rate function). Suppose that an infinite line, say  $\{z_0 + tz\}_{t \in \mathbb{R}}$  is contained in  $\partial\Lambda(\eta)$ . Then, for all  $\theta \in \mathbb{R}^d$  and any  $t \in \mathbb{R}$ ,

$$\Lambda(\theta) \geq \Lambda(\eta) + \langle \theta - \eta, z_0 \rangle + t \langle \theta - \eta, z \rangle \quad (2.3.99)$$

which is possible only if  $\mathcal{D}_\Lambda \subseteq \{\theta : \langle \theta, z \rangle = 0\}$  contradicting Assumption 2.3.1. From the above, it follows that  $\partial\Lambda(\eta)$  is a convex closed set which does not contain infinite lines and as such by [24],

18.5.3, it contains an exposed point, i.e. there exists  $x \in \partial\Lambda(\eta)$  and a vector  $\mathbf{v} \in \mathbb{R}^d$  such that

$$\langle \mathbf{v}, x \rangle > \langle \mathbf{v}, z \rangle \quad \forall z \in \partial\Lambda(\eta), \quad z \neq x. \quad (2.3.100)$$

Let the normal cone to  $\mathcal{D}_\Lambda$  at  $\eta$  be defined as

$$\mathcal{N} \triangleq \{\mathbf{n} : \langle \lambda - \eta, \mathbf{n} \rangle \leq 0 \quad \text{for all } \lambda \in \mathcal{D}_\Lambda\} \quad (2.3.101)$$

and note that  $\mathcal{N}$  is non-empty since  $\mathcal{D}_\Lambda$  is a convex set of non-empty interior and  $\eta \in \mathcal{D}_\Lambda \setminus \mathcal{D}_\Lambda^\circ$ . Then, for any  $\mathbf{n} \in \mathcal{N}$  and any  $\lambda \in \mathbb{R}^d$

$$\Lambda(\lambda) \geq \Lambda(\eta) + \langle \lambda - \eta, \mathbf{n} \rangle \geq \Lambda(\eta) + \langle \lambda - \eta, x + \mathbf{n} \rangle,$$

so that  $x + \mathbf{n} \in \partial\Lambda(\eta)$  where  $x \in \partial\Lambda(\eta)$  is the exposed point defined above. Therefore, in particular, by (2.3.100)

$$\langle \mathbf{v}, \mathbf{n} \rangle < 0 \quad \forall \mathbf{n} \in \mathcal{N}. \quad (2.3.102)$$

We next claim that  $\eta + \delta\mathbf{v} \in \mathcal{D}_\Lambda^\circ$  for all  $\delta > 0$  small enough. Indeed, assume otherwise, then by [24], 23.7.1 there exists  $\mathbf{n} \in \mathcal{N}$  such that

$$\sup_{\lambda \in \mathcal{D}_\Lambda} \langle \lambda, \mathbf{n} \rangle \leq \langle \eta, \mathbf{n} \rangle \leq \langle \eta + \delta\mathbf{v}, \mathbf{n} \rangle \quad (2.3.103)$$

contradicting the inequality (2.3.102) above.

Choose now a sequence  $\lambda_n = \eta + \delta_n\mathbf{v} \in \mathcal{D}_\Lambda^\circ$  such that  $\delta_n \rightarrow 0$ . Since  $\nabla\Lambda(\lambda_n) \in \partial\Lambda(\lambda_n)$  for any  $n$  it follows that

$$\Lambda(\theta) \geq \Lambda(\lambda_n) + \langle \theta - \lambda_n, \nabla\Lambda(\lambda_n) \rangle, \quad \forall \theta \in \mathbb{R}^d \quad (2.3.104)$$

Because of the assumption that  $\Lambda$  is a steep function,  $\epsilon_n \triangleq 1/|\nabla\Lambda(\lambda_n)| \rightarrow 0$  as  $n \rightarrow \infty$  and  $\epsilon_n \nabla\Lambda(\lambda_n)$  has a limit point  $y \in \mathbb{R}^d$  with  $|y| = 1$ . Passing to the convergent subsequence  $\{\lambda_n\}$ , for any  $n$  large enough and any  $\theta \in \mathbb{R}^d$ ,

$$\begin{aligned} \Lambda(\theta) &\geq (1 - \epsilon_n) \Lambda(\eta) + (1 - \epsilon_n) \langle \theta - \eta, x \rangle \\ &\quad + \epsilon_n \Lambda(\lambda_n) + \langle \theta - \lambda_n, \epsilon_n \nabla\Lambda(\lambda_n) \rangle \end{aligned} \quad (2.3.105)$$

where  $x \in \partial\Lambda(\eta)$  is as specified in (2.3.100). In the limit  $n \rightarrow \infty$  by the upper semicontinuity of the convex function  $\Lambda(\cdot)$

$$\Lambda(\theta) \geq \Lambda(\eta) + \langle \theta - \eta, x + y \rangle \quad (2.3.106)$$

(as  $\limsup_{n \rightarrow \infty} \Lambda(\lambda_n) \leq \Lambda(\eta) < \infty$ ). Therefore, in particular

$$\langle v, y \rangle < 0 \quad (2.3.107)$$

By comparing (2.3.104) for  $\theta = \lambda_n = \eta + \delta_n v$  with

$$\Lambda(\lambda_n) \geq \Lambda(\eta) + \langle \lambda_n - \eta, x \rangle = \Lambda(\eta) + \delta_n \langle v, x \rangle$$

one obtains

$$\delta_n \langle v, x \rangle \leq \delta_n \langle v, \nabla \Lambda(\lambda_n) \rangle \quad (2.3.108)$$

for the same subsequence  $\{\lambda_n\}$  as used above. Multiplying by  $\frac{\epsilon_n}{\delta_n} > 0$  and taking the limit  $n \rightarrow \infty$  yields  $\langle v, y \rangle \geq 0$  in contradiction with (2.3.107). In conclusion,  $\partial \Lambda(\eta) = \emptyset$  for any  $\eta \in \mathcal{D}_\Lambda \setminus \mathcal{D}_\Lambda^\circ$ .  $\square$

### Exercises:

**2.3.1** Prove that for any  $\mu \in M_1(\Sigma)$ , the relative entropy  $H(\cdot|\mu)$  is the Fenchel-Legendre transform of the function  $\Lambda(\cdot)$  defined in (2.3.82).

**Hint:** Prove first that  $\mathcal{D}_{\Lambda^*} = M_1(\Sigma_\mu)$ . Then, show that  $\nu(a_i) = 0$  and  $\nu \in M_1(\Sigma_\mu)$  imply that the value of  $\Lambda^*(\nu)$  is obtained by taking  $\lambda_i = -\infty$ . Finally, show that  $\nu = \nabla \Lambda(\eta)$  when  $\nu$  is a probability vector with  $\Sigma_\nu = \Sigma_\mu = \Sigma$  and

$$\eta_i \triangleq \log \left[ \frac{\nu(a_i)}{\mu(a_i)} \right] \quad i = 1, \dots, |\Sigma| \quad (2.3.109)$$

Conclude that then  $\Lambda^*(\nu) = \langle \eta, \nu \rangle - \Lambda(\eta) = H(\nu|\mu)$ .

**2.3.2 (a).** Use the exponential form of Markov's inequality to prove that for any  $C \subset \mathbb{R}^d$  any  $n$  and any  $\lambda \in \mathbb{R}^d$

$$a_n \log \mu_n(C) \leq - \inf_{y \in C} \langle \lambda, y \rangle + a_n \Lambda_n(a_n^{-1} \lambda).$$

(b). Assume that  $\Lambda(\lambda) = a_n \Lambda_n(a_n^{-1} \lambda)$  for any  $n$  (examples where this is true are given in Theorem 2.3.2 and in exercise 2.3.5). Recall the following version of the min-max theorem: let  $g(\theta, y)$  be convex and lower semicontinuous in  $y$ , concave and upper semicontinuous in  $\theta$ . Let  $C \subset \mathbb{R}^d$  be convex and compact. Then

$$\inf_{y \in C} \sup_{\theta} g(\theta, y) = \sup_{\theta} \inf_{y \in C} g(\theta, y)$$



(c.f. [10], pg. 174). Apply this theorem to justify the upper bound

$$a_n \log \mu_n(C) \leq - \sup_{\lambda \in \mathbb{R}^d} \inf_{y \in C} [\langle \lambda, y \rangle - \Lambda(\lambda)] = - \inf_{y \in C} \Lambda^*(y)$$

for any  $n$  and any convex, compact set  $C$ .

(c). Establish (2.3.64) for all compact sets by applying the above bound. Note that this approach yields a concrete upper bound for any finite  $n$ .

(d). Find a compact set  $B$  for which

$$\sup_{\lambda \in \mathbb{R}^d} \inf_{y \in B} [\langle \lambda, y \rangle - \Lambda(\lambda)] < \inf_{y \in B} \Lambda^*(y).$$

**2.3.3** Prove that in the assumptions of Theorem 2.3.2, the assumption of lower semicontinuity of  $\Lambda$  may be dropped as it follows from the steepness and i.i.d. structure.

**Hint** you need to extend exercise 2.2.4 to  $\mathbb{R}^d$ . The only difficulty is with sequences  $\lambda_n \rightarrow \lambda$  not along a line, but such that their angle from a line converges to zero. Use the last observation to bound the difference between these two situations in terms of the function  $G$  introduced in exercise 2.2.4.

**2.3.4** Let  $(w_{t_1}, \dots, w_{t_d})$  be samples of a Brownian motion at  $t_1, \dots, t_d$ , i.e.  $(w_{t_{j+1}} - w_{t_j})$  is a Normal random variable independent of  $w_{t_\ell}$ ,  $\ell \leq j$ , of variance  $(t_{j+1} - t_j)$  and zero mean. Find the rate function for the empirical mean  $\hat{S}_n$  of  $X_i \triangleq (w_{t_1}^i, \dots, w_{t_d}^i)$  where  $w_{t_j}^i$ ,  $i = 1, \dots, n$  are samples of independent Brownian motions at instances  $t_j$ .

**Remark:** Note that the law of  $\hat{S}_n$  is the same as that of  $\frac{1}{\sqrt{n}}(w_{t_1}, \dots, w_{t_d})$ , and compare to Schilder's Theorem which is presented in Section ??.

**2.3.5** Let  $X_j$  be i.i.d. random variables over  $\mathbb{R}^d$  with a steep logarithmic moment generating function  $\Lambda$  such that  $0 \in \mathcal{D}_\Lambda^\circ$ . Let  $N(t)$  be a Poisson process of unit rate which is independent of the  $X_j$  variables and consider the random variables

$$\hat{S}_n \triangleq \frac{1}{n} \sum_{j=1}^{N(n)} X_j.$$

Prove that the family of laws  $\mu_n$  corresponding to  $\hat{S}_n$  satisfies a large deviations principle with the rate function being the Fenchel-Legendre transform of  $e^{\Lambda(\lambda)} - 1$ .

Hint: You can apply Theorem 2.3.2 as  $N(n) = \sum_{j=1}^n N_j$  where  $N_j$  are i.i.d. Poisson(1) random variables.

**2.3.6** Let  $N(n)$  be a sequence of integer valued random variables whose logarithmic moment generating functions  $\Lambda_n$  satisfy the Assumption 2.3.1. Let  $X_j$  be i.i.d. random variables over  $\mathbb{R}^d$  with *finite everywhere* logarithmic moment generating function  $\Lambda_X$  and let  $\mu_n$  denotes the law of

$$Z_n \triangleq a_n \sum_{j=1}^{N(n)} X_j .$$

Prove that if the conditions of Lemma 2.3.2 hold for the convex function  $\Lambda(\Lambda_X(\lambda))$  then  $\mu_n$  satisfies a large deviations principle governed by the Fenchel-Legendre transform of this function.

**2.3.7** For any  $\delta > 0$  let  $Z_{n,\delta} = Z_n + \sqrt{\delta a_n} V$  where  $V$  is a standard multivariate Normal random variable.

(a). Prove that when Assumption 2.3.1 holds for  $Z_n$  it also holds for  $Z_{n,\delta}$  with the limiting logarithmic moment generating function  $\Lambda_\delta(\lambda) = \Lambda(\lambda) + \frac{\delta}{2} |\lambda|^2$ .

(b). Show that for any  $x \in \mathbb{R}^d$  the value of the Fenchel-Legendre transform of  $\Lambda_\delta$  does not exceed  $\Lambda^*(x)$ .

(c). Prove that if  $\lim_{n \rightarrow \infty} E(Z_n) = \bar{x} \in \mathbb{R}^d$  exists then  $\Lambda(\lambda) \geq \langle \lambda, \bar{x} \rangle$  for any  $\lambda \in \mathbb{R}^d$ . Conclude that if in addition  $\Lambda$  is finite and differentiable everywhere then  $\mathcal{F}_\delta = \mathbb{R}^d$  (where  $\mathcal{F}_\delta \triangleq \{x : x = \nabla \Lambda_\delta(\lambda) \text{ for some } \lambda \in \mathbb{R}^d\}$ ).

(d). By applying Theorem 2.3.1 for  $Z_{n,\delta}$  and (a)-(c) above deduce that for any  $x \in \mathbb{R}^d$  and any  $\epsilon > 0$

$$\liminf_{n \rightarrow \infty} a_n \log \text{Prob}(Z_{n,\delta} \in B_{x,\epsilon/2}) \geq -\Lambda^*(x) \quad (2.3.110)$$

(e). Prove that

$$\limsup_{n \rightarrow \infty} a_n \log \text{Prob}(\sqrt{\delta a_n} |V| \geq \epsilon/2) \leq -\frac{\epsilon^2}{8\delta} \quad (2.3.111)$$

(f). Prove that

$$\text{Prob}(Z_n \in B_{x,\epsilon}) \geq \text{Prob}(Z_{n,\delta} \in B_{x,\epsilon/2}) - \text{Prob}(\sqrt{\delta a_n} |V| \geq \epsilon/2) \quad (2.3.112)$$

and by combining (2.3.110), (2.3.111) and (2.3.112) (for  $n \rightarrow \infty$  and then  $\delta \rightarrow 0$ ) conclude that the large deviations lower bound holds for the laws  $\mu_n$  corresponding to  $Z_n$ .

(g). Deduce now by part (a) of Theorem 2.3.1 that when Assumption 2.3.1 holds with  $\Lambda$  which is finite and differentiable everywhere and when moreover  $\lim_{n \rightarrow \infty} E(Z_n)$  exists then  $\mu_n$  satisfy a large deviations principle with rate function  $\Lambda^*$ .

**Remark:** This may serve for example as an alternative derivation of Cramer's theorem which avoids the convex analysis Lemma 2.3.2.

**2.3.8** Let  $X_1, \dots, X_n, \dots$  be a real-valued, zero mean, stationary Gaussian process with covariance sequence  $R_i \triangleq E(X_n X_{n+i})$ . Suppose the process has a finite power  $P$  defined via  $P \triangleq \lim_{n \rightarrow \infty} \sum_{i=-n}^n R_i (1 - \frac{|i|}{n})$ . Let  $\mu_n$  be the law of the empirical mean  $\hat{S}_n$  of the first  $n$  samples of this process. Prove that  $\mu_n$  satisfy a large deviations principle controlled by the good rate function  $\Lambda^*(x) = \frac{x^2}{2P}$ .

**2.3.9** Again, let  $X_1, \dots, X_n, \dots$  be a real-valued, zero mean, stationary Gaussian process with covariance sequence  $R_i \triangleq E(X_n X_{n+i})$ . Assume that this covariance sequence is absolutely summable and let  $S(\omega) > 0$  denote its Fourier transform. Consider the empirical covariances

$$Z_n^j = \frac{1}{n} \sum_{i=1}^{n-j} X_i X_{i+j}$$

for  $j = 0, \dots, d-1$ . Let  $Z_n \in \mathbb{R}^d$  be the empirical covariance vector composed of  $\{Z_n^j\}$  and  $\Lambda_n$  the corresponding logarithmic moment generating functions.

(a). Verify that

$$\Lambda_n(n\theta) = -\frac{1}{2} \sum_{i=1}^n \log[1 - \lambda_i(\Theta \mathbf{R})].$$

Here  $\lambda_i(\Theta \mathbf{R})$  is the  $i$ -th eigenvalue of the product of the covariance matrix  $\mathbf{R}$  and the matrix  $\Theta$  where  $\Theta(j, k) = \theta_{j-k}$  for all  $j \in \{k, \dots, k+d-1\}$  and is zero otherwise.

(b). Assume that

$$\lim_{n \rightarrow \infty} -\frac{1}{2n} \sum_{i=1}^n \log[1 - \lambda_i(\Theta \mathbf{R})] = -\frac{1}{4\pi} \int_0^{2\pi} \log[1 - S(\omega) \sum_{k=0}^{d-1} \theta_k e^{-i\omega k}] d\omega \triangleq \Lambda(\theta)$$

(this identity indeed holds by the limiting distribution of near Toeplitz matrices, see [16]). Prove that the empirical covariance vectors  $Z_n$  satisfy a large deviations principle controlled by the Fenchel-Legendre transform of the function  $\Lambda$  defined above.

## 2.4 Large deviations of Markov chains over finite alphabets

The results of Section 2.1 are extended in this section to random variables  $X_1, \dots, X_n$  which take values in the finite alphabet  $\Sigma = \{a_1, \dots, a_{|\Sigma|}\}$ , with a Markov structure instead of an i.i.d. structure. Although most of the results may be derived by the method of types presented in Section 2.1, the combinatorics involved are quite cumbersome (for more details about this approach consider exercise 2.4.7). Thus, an alternative derivation of these results via an application of Theorem 2.3.1 is adopted here. Without loss of generality, identify  $\Sigma$  with the set  $\{1, \dots, |\Sigma|\}$  so that  $a_i = i$ .

Let  $\Pi = \{\pi(i, j)\}_{i, j=1, \dots, |\Sigma|}$  be a stochastic matrix, i.e. a matrix whose elements are non-negative and such that each row-sum is one.  $P_x^\pi$  denotes the Markov probability measure associated with the transition probability  $\Pi$  and initial state  $x \in \Sigma$ . Specifically,

$$P_x^\pi(X_1 = x_1, \dots, X_n = x_n) = \pi(x, x_1) \prod_{i=1}^{n-1} \pi(x_i, x_{i+1}). \quad (2.4.113)$$

A matrix  $B$  with nonnegative entries is called *irreducible*, if for any pair of indices  $i, j$  there exists an  $m(i, j)$  such that  $B^{m(i, j)}(i, j) > 0$ , where  $B^m$  denotes the usual product of matrices. This property is equivalent to the condition that one may find for each  $i, j$  a sequence of indices  $i_1, \dots, i_m$  such that  $i_1 = i$ ,  $i_m = j$  and  $B(i_k, i_{k+1}) > 0$  for all  $k = 1, \dots, m-1$ . The following theorem describes basic properties of irreducible matrices.

**Theorem 2.4.1 (Perron-Frobenius)** *Let  $B = \{B(i, j)\}_{i, j=1}^{|\Sigma|}$  be an irreducible matrix. Then there exists an eigenvalue  $\rho$  (called the Perron-Frobenius eigenvalue) such that*

(a).  $\rho > 0$  is real.

(b). *There exist right and left eigenvectors corresponding to the eigenvalue  $\rho$  which are strictly positive, i.e. there exist vectors  $\mu, \vartheta$  with  $\mu_i, \vartheta_i > 0$  for all  $i$  (this is denoted in the sequel by  $\mu \gg 0, \vartheta \gg 0$ ) such that*

$$\sum_{j=1}^{|\Sigma|} B(i, j) \vartheta_j = \rho \vartheta_i \quad (2.4.114)$$

$$\sum_{i=1}^{|\Sigma|} \mu_i B(i, j) = \rho \mu_j \quad (2.4.115)$$

(c). For any eigenvalue  $\lambda$  of  $B$ ,  $|\lambda| \leq \rho$ .

(d). The right and left eigenvectors  $\mu, \nu$  corresponding to the eigenvalue  $\rho$  are unique up to a constant multiple.

(e). Let  $\phi$  be any strictly positive vector, then for any  $i \in \Sigma$

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \left[ \sum_{j=1}^{|\Sigma|} B^n(i, j) \phi_j \right] = \lim_{n \rightarrow \infty} \frac{1}{n} \log \left[ \sum_{j=1}^{|\Sigma|} B^n(j, i) \phi_j \right] = \log \rho \quad (2.4.116)$$

**Proof:** Parts (a)÷(d) are stated for example in [27], Theorem 1.5. To prove part (e), let  $\alpha \triangleq \sup_i \nu_i$ ,  $\beta \triangleq \inf_i \nu_i > 0$  and  $\gamma \triangleq \inf_j \phi_j > 0$ ,  $\delta = \sup_j \phi_j$  (where  $\nu$  is the right eigenvector corresponding to  $\rho$  above). Then, for all  $i, j \in \Sigma$ ,

$$\frac{\delta}{\beta} B^n(i, j) \nu_j \geq B^n(i, j) \phi_j \geq \frac{\gamma}{\alpha} B^n(i, j) \nu_j \quad (2.4.117)$$

Therefore,

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{n} \log \left[ \sum_{j=1}^{|\Sigma|} B^n(i, j) \phi_j \right] &= \lim_{n \rightarrow \infty} \frac{1}{n} \log \left[ \sum_{j=1}^{|\Sigma|} B^n(i, j) \nu_j \right] \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} \log (\rho^n \nu_i) = \log \rho. \end{aligned} \quad (2.4.118)$$

A similar argument leads to

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \left[ \sum_{j=1}^{|\Sigma|} B^n(j, i) \phi_j \right] = \log \rho. \quad (2.4.119)$$

□

### 2.4.1 Cramer's theorem for Markov additive processes

Let  $g : \Sigma \rightarrow \mathbb{R}^d$  be a given deterministic function. The large deviations of the empirical  $g$ -mean

$$Z_n = \frac{1}{n} \sum_{k=1}^n g(X_k) \quad (2.4.120)$$

are the subject of this section (for the extension to random functions see exercise 2.4.1). If  $X_k$  were independent, Cramer's Theorem 2.1.1 applies and then the large deviations principle is governed

by the Legendre transform of the logarithmic moment generating function. Theorem 2.3.1 hints that the rate function may still be expressed in terms of a Legendre transform even in the current dependent case where the  $X_k$  possess a Markov structure.

In order to find the proper function which replaces the logarithmic moment generating function  $\Lambda(\lambda)$  associate with any  $\lambda \in \mathbb{R}^d$  a non-negative matrix  $\Pi_\lambda$  via

$$\pi_\lambda(i, j) = \pi(i, j) e^{<\lambda, g(j)>} \quad i, j \in \Sigma. \quad (2.4.121)$$

When  $\Pi$  is irreducible it follows that  $\Pi_\lambda$  are also irreducible matrices since  $e^{<\lambda, g(j)>}$  is positive. Let  $\rho(\Pi_\lambda)$  denotes the Perron-Frobenius eigenvalue of  $\Pi_\lambda$  then  $\log \rho(\Pi_\lambda)$  plays the role of the logarithmic moment generating function  $\Lambda(\lambda)$ . Specifically, the following analog of Theorem 2.3.2 holds.

**Theorem 2.4.2** *Assume  $\Pi$  is irreducible and define*

$$I(z) \triangleq \sup_{\lambda \in \mathbb{R}^d} \{ <\lambda, z> - \log \rho(\Pi_\lambda) \} \quad (2.4.122)$$

*Then,  $I(\cdot)$  is a good, convex, rate function which controls the large deviations of the empirical  $g$ -means  $\{Z_n\}$ , i.e. for any measurable set  $\Gamma \subseteq \mathbb{R}^d$ , and any initial state  $x \in \Sigma$ ,*

$$- \inf_{z \in \Gamma^o} I(z) \leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log P_x^\pi(\{Z_n \in \Gamma\}) \leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log P_x^\pi(\{Z_n \in \Gamma\}) \leq - \inf_{z \in \Gamma} I(z) \quad (2.4.123)$$

**Proof:** Define

$$\Lambda_n(\lambda) \triangleq \log E_x^\pi \left[ e^{<\lambda, Z_n>} \right] \quad (2.4.124)$$

In view of Theorem 2.3.1, it is enough to check that the limit

$$\Lambda(\lambda) \triangleq \lim_{n \rightarrow \infty} \frac{1}{n} \Lambda_n(n\lambda) = \lim_{n \rightarrow \infty} \frac{1}{n} \log E_x^\pi \left[ e^{n<\lambda, Z_n>} \right] \quad (2.4.125)$$

exists for all  $\lambda \in \mathbb{R}^d$ , is differentiable, and that  $\Lambda(\lambda) = \log \rho(\Pi_\lambda)$ . Note that

$$\begin{aligned} \Lambda_n(n\lambda) &= \log E_x^\pi \left[ e^{<n\lambda, \sum_{k=1}^n g(X_k)>} \right] \\ &= \log \sum_{x_1, \dots, x_n} P_x^\pi(X_1 = x_1, \dots, X_n = x_n) \prod_{k=1}^n e^{<\lambda, g(x_k)>} \\ &= \log \sum_{x_1, \dots, x_n} \pi(x, x_1) e^{<\lambda, g(x_1)>} \dots \pi(x_{n-1}, x_n) e^{<\lambda, g(x_n)>} \\ &= \log \sum_{j=1}^{|\Sigma|} (\Pi_\lambda)^n(x, j) \end{aligned} \quad (2.4.126)$$

Since  $\Pi_\lambda$  is an irreducible matrix, part (e) of the Perron-Frobenius theorem yields (for  $\phi_j = 1$ )

$$\Lambda(\lambda) = \lim_{n \rightarrow \infty} \frac{1}{n} \Lambda_n(n\lambda) = \log \rho(\Pi_\lambda) \quad (2.4.127)$$

Moreover, since  $|\Sigma|$  is finite it is clear that  $\rho(\Pi_\lambda)$ , being an isolated root of the characteristic equation for the matrix  $\Pi_\lambda$  is differentiable with respect to  $\lambda$  (see [14] for details). Therefore, Theorem 2.3.1 may be applied to complete the proof.  $\square$

**Remark:** The above proof relies on two properties of the Markov chain – namely, part (e) of the Perron-Frobenius theorem and the differentiability of  $\rho(\Pi_\lambda)$  with respect to  $\lambda$ . Thus, Theorem 2.4.2 holds as long as the Markov chain has these two properties. In particular, the finiteness of  $\Sigma$  is not crucial and indeed a large deviations principle for a general Markov chain set-up is presented in Section ??.

**Exercises:**

2.4.1 Assume that  $X_1, \dots, X_n$  have the joint law  $P_x^\pi$  where  $\Pi$  is an irreducible stochastic matrix. Consider the empirical means

$$Z_n = \frac{1}{n} \sum_{k=1}^n Y_k$$

where the conditional law of  $Y_k$  when  $X_k = j$  is  $\mu_j \in M_1(\mathbb{R}^d)$  and for any given realization  $\{X_k\}_{k=1}^n$  of the Markov chain states the variables  $Y_k$  are conditionally independent. Suppose that the logarithmic moment generating functions  $\Lambda_j$  associated with  $\mu_j$  are finite everywhere (for all  $j \in \Sigma$ ). Prove that Theorem 2.4.2 holds in this case where now

$$\pi_\lambda(i, j) \triangleq \pi(i, j) e^{\Lambda_j(\lambda)} \quad i, j \in \Sigma .$$

## 2.4.2 Sanov's theorem for the empirical measure of Markov chains

A particularly important application of Theorem 2.4.2 above yields the large deviations principle satisfied by the empirical measure of Markov chains. Namely, define  $L_n^{\mathbf{X}}(i) = \frac{1}{n} \sum_{j=1}^n \delta_i(X_j)$ , where

$$\delta_i(x) \triangleq \begin{cases} 1 & i = x \\ 0 & \text{otherwise} \end{cases} .$$

For  $i = 1, \dots, |\Sigma|$ , let  $a_n(i) = E_x^\pi[L_n^{\mathbf{X}}(i)] = \frac{1}{n} \sum_{j=1}^n \pi^j(x, i)$ . Then,  $a_n \rightarrow \mu$  as  $n \rightarrow \infty$  where  $\mu$  is a unique properly normalized left eigenvector of  $\Pi$  (since  $|a_n(\Pi - \mathbf{I})(j)| = \frac{1}{n} |\pi^{n+1}(x, j) - \pi(x, j)| \leq \frac{2}{n}$ ).

Further, by Chebychev's inequality,  $L_n^{\mathbf{X}} \rightarrow \mu$  as  $n \rightarrow \infty$  and therefore,  $L_n^{\mathbf{X}}$  is a good candidate for a large deviations statement on  $M_1(\Sigma)$ .

It is clear that  $L_n^{\mathbf{X}}$  fits into the framework of Section 2.4.1 with  $g(x) = (\delta_1(x), \dots, \delta_{|\Sigma|}(x))$ . Therefore, by Theorem 2.4.2, a large deviations principle follows with the rate function

$$I(q) = \sup_{\lambda \in \mathbb{R}^d} (\langle \lambda, q \rangle - \log \rho(\Pi_\lambda)), \quad (2.4.128)$$

where here  $\pi_\lambda(i, j) \triangleq \pi(i, j) e^{\lambda_j}$  and  $q \in M_1(\Sigma)$ . The following alternative characterization of  $I(q)$  is sometimes more useful.

**Theorem 2.4.3**

$$I(q) = J(q) \triangleq \begin{cases} \infty & q \notin M_1(\Sigma) \\ \sup_{u \gg 0} \sum_{j=1}^{|\Sigma|} q_j \log \left[ \frac{u_j}{(u\Pi)_j} \right] & q \in M_1(\Sigma) \end{cases} \quad (2.4.129)$$

**Remarks:** This identity actually holds also for non-stochastic matrices (see exercise 2.4.3). In the i.i.d. set up the rows of  $\Pi$  are identical and then  $J(q)$  is merely the relative entropy  $H(q|\pi(1, \cdot))$  (see exercise 2.4.2).

**Proof:** Since  $M_1(\Sigma)$  is a closed subset of  $\mathbb{R}^{|\Sigma|}$  its complement  $M_1(\Sigma)^c$  is an open set. Further,  $L_n^{\mathbf{X}} \in M_1(\Sigma)$ , for any  $n$  and any realization  $\mathbf{X}$ . Therefore, by the large deviations lower bound (2.4.123)

$$-\infty = \frac{1}{n} \log P_x^\pi(\{L_n^{\mathbf{X}} \in M_1(\Sigma)^c\}) \geq - \inf_{q \notin M_1(\Sigma)} I(q), \quad (2.4.130)$$

i.e.,  $I(q) = \infty$  for any  $q \notin M_1(\Sigma)$ .

Fix  $q \in M_1(\Sigma)$ ,  $u \gg 0$  and set  $\lambda_j = \log \left[ \frac{u_j}{(u\Pi)_j} \right]$  (since  $u \gg 0$  and  $\Pi$  is irreducible, it follows that  $u\Pi \gg 0$ ). Observe that  $u\Pi_\lambda = u$  and thus  $\rho(\Pi_\lambda) = 1$  by part (e) of the Perron-Frobenius theorem (with  $\phi_i = u_i > 0$ ). Therefore,

$$I(q) \geq \sum_{j=1}^{|\Sigma|} q_j \log \left[ \frac{u_j}{(u\Pi)_j} \right] - \log 1$$

Since  $u \gg 0$  is arbitrary this inequality implies that  $I(q) \geq J(q)$ .



To establish the reverse inequality, fix an arbitrary vector  $\lambda \in \mathbb{R}^{|\Sigma|}$  and let  $\alpha > \rho(\Pi_\lambda)$  be any upper bound on  $\rho(\Pi_\lambda)$ . Define

$$u_j^* \triangleq \sum_{n=0}^{\infty} \alpha^{-n} \sum_{i=1}^{|\Sigma|} (\Pi_\lambda)^n(i, j) < \infty, \quad j = 1, \dots, |\Sigma|, \quad (2.4.131)$$

where the finiteness of  $u_j^*$  is a direct consequence of part (e) of the Perron-Frobenius theorem. Moreover,  $u^* \Pi_\lambda = \alpha(u^* - 1)$  (where  $1$  denotes the all ones vector) and  $u^* \gg 0$ . Thus, by the definition of  $\Pi_\lambda$ ,

$$\begin{aligned} \langle \lambda, q \rangle &+ \sum_{j=1}^{|\Sigma|} q_j \log \frac{(u^* \Pi)_j}{u_j^*} = \sum_{j=1}^{|\Sigma|} q_j \log \frac{(u^* \Pi_\lambda)_j}{u_j^*} \\ &= \sum_{j=1}^{|\Sigma|} q_j \log \frac{\alpha(u^* - 1)_j}{u_j^*} = \log \alpha + \sum_{j=1}^{|\Sigma|} q_j \log \left( 1 - \frac{1}{u_j^*} \right) < \log \alpha \end{aligned} \quad (2.4.132)$$

Thus,  $(\langle \lambda, q \rangle - \log \alpha) < \sum_{j=1}^{|\Sigma|} q_j \log \frac{u_j^*}{(u^* \Pi)_j}$ . Take  $\alpha \downarrow \rho(\Pi_\lambda)$  to deduce that

$$\langle \lambda, q \rangle - \log \rho(\Pi_\lambda) \leq \sup_{u \gg 0} \sum_{j=1}^{|\Sigma|} q_j \log \frac{u_j}{(u \Pi)_j} = J(q) \quad (2.4.133)$$

and since  $\lambda$  is arbitrary, also  $I(q) \leq J(q)$  and the proof is complete.  $\square$

#### Exercises:

**2.4.2** Suppose  $\pi(i, j) = \mu(j)$ ,  $i, j \in \Sigma$  where  $\mu \in M_1(\Sigma)$ . Prove that then  $J(\cdot) = H(\cdot | \mu)$  (the relative entropy with respect to  $\mu$ ) while  $I(\cdot)$  is the Fenchel-Legendre transform of the moment generating function  $\Lambda$  of (2.3.82). Thus, Theorem 2.4.3 is the natural extension of exercise 2.3.1 to the Markov set-up.

**2.4.3 (a).** Show that (2.4.129) holds for any non-negative irreducible matrix  $\Pi$  (not necessarily stochastic).

**Hint:** Let  $\phi(i) = \sum_j \pi(i, j)$ . Clearly,  $\phi \gg 0$ , and thus the matrix  $\Pi^*$  given by,  $\pi^*(i, j) = \pi(i, j)/\phi(i)$  is stochastic. Prove now that  $J_{\Pi^*}(q) = J_\Pi(q) + \sum_j q_j \log \phi(j)$  for any  $q \in \mathbb{R}^{|\Sigma|}$ , and likewise  $I_{\Pi^*}(q) = I_\Pi(q) + \sum_j q_j \log \phi(j)$  (where  $J_\Pi$  and  $I_\Pi$  denote the rate functions  $J$  and  $I$  associated with the matrix  $\Pi$  via (2.4.129) and (2.4.128) respectively).

(b). Show that for any irreducible, non-negative matrix  $\Pi$

$$\log \rho(\Pi) = \sup_{\nu \in M_1(\Sigma)} \{-J(\nu)\} \quad (2.4.134)$$

This characterization of the spectral radius of non-negative matrices is useful when looking for tight bounds (for an alternative characterization see exercise 2.4.5).

### 2.4.3 Sanov's theorem for the pair empirical measure of Markov chains

The large deviations principle for the empirical measure of a Markov chain is still in the form of an optimization problem. Moreover the nice interpretation in terms of entropy (recall Section 2.1.1 where the i.i.d. case is presented) has disappeared. It is interesting to note that by considering a somewhat different random variable, from which the large deviations for  $L_n^X$  may be recovered (see exercise 2.4.4), one is also able to get a large deviations with a rate function which is an appropriate relative entropy.

Consider the space  $\Sigma^{(2)} \triangleq \Sigma \times \Sigma$ , which corresponds to consecutive pairs of elements from the sequence  $X$ . Note that by considering the pairs formed by  $X_1, \dots, X_n$ , i.e. the sequence  $X_1X_2, X_2X_3, \dots, X_iX_{i+1}, \dots, X_{n-1}X_n$ , one recovers a Markov chain with state space  $\Sigma^{(2)}$  and transition matrix  $\Pi^{(2)}$  specified via

$$\pi^{(2)}(k \times \ell, i \times j) = \delta_\ell(i) \pi(i, j) \quad (2.4.135)$$

For simplicity assume throughout this section that  $\Pi$  is strictly positive (i.e.  $\pi(i, j) > 0$  for all  $i, j$ ). Then,  $\Pi^{(2)}$  is an irreducible transition matrix, and therefore the results of Section 2.4.2 may be applied to find the large deviations rate function  $I^{(2)}(q)$  associated with the pair empirical measures

$$L_n^{X, (2)}(y) \triangleq \frac{1}{n} \sum_{i=1}^n \delta_y(X_{i-1}, X_i), \quad y \in \Sigma^{(2)}. \quad (2.4.136)$$

Note that  $L_n^{X, (2)} \in M_1(\Sigma^{(2)})$  and therefore  $I^{(2)}(\cdot)$  is a good, convex, rate function over this space. The next theorem characterizes  $I^{(2)}(\cdot)$  as an appropriate relative entropy. The following definitions are needed for that purpose. For any  $q \in M_1(\Sigma^{(2)})$ , let  $q(i) \triangleq \sum_{j=1}^{|\Sigma|} q(i, j)$  be its marginal and when  $q(i) > 0$  let  $q(j|i) \triangleq \frac{q(i, j)}{q(i)}$ . A measure  $q \in M_1(\Sigma^{(2)})$  is called **shift invariant** if  $q(i) = \sum_{k=1}^{|\Sigma|} q(k, i)$  for all  $i$  (i.e., both marginals of  $q$  are identical).

**Theorem 2.4.4** *Assume  $\Pi$  is strictly positive. Then for any  $q \in M_1(\Sigma^{(2)})$ ,*

$$I^{(2)}(q) = \begin{cases} \sum_i q(i) H(q(\cdot|i) | \pi(i, \cdot)), & q \text{ shift invariant} \\ \infty & \text{otherwise} \end{cases} \quad (2.4.137)$$

where  $H(\cdot|\cdot)$  is the relative entropy function defined in Section 2.1.1. Specifically,

$$H(q(\cdot|i)|\pi(i, \cdot)) = \sum_{j=1}^{|\Sigma|} q(j|i) \log \frac{q(j|i)}{\pi(i, j)}. \quad (2.4.138)$$

**Remarks:** When  $\Pi$  is not strictly positive (but is irreducible) the theorem still applies with  $\Sigma^{(2)}$  replaced by  $\Sigma_{\Pi} \triangleq \{(i, j) : \pi(i, j) > 0\}$ , and an almost identical proof. The above representation of  $I^{(2)}(q)$  is useful for example in characterizing the spectral radius of non-negative matrices (see exercise 2.4.5) and in establishing the analog of Sanov's theorem for *time weighted* empirical measures (see exercise 2.4.6). It is also useful because bounds on the relative entropy are readily available and may be used to obtain bounds on the rate function.

**Proof:** By Theorem 2.4.3

$$\begin{aligned} I^{(2)}(q) &= \sup_{u \gg 0} \sum_{j=1}^{|\Sigma|} \sum_{i=1}^{|\Sigma|} q(i, j) \log \frac{u(i, j)}{(u\Pi^{(2)})(i, j)} \\ &= \sup_{u \gg 0} \sum_{j=1}^{|\Sigma|} \sum_{i=1}^{|\Sigma|} q(i, j) \log \frac{u(i, j)}{\left[ \sum_k u(k, i) \right] \pi(i, j)} \end{aligned} \quad (2.4.139)$$

where the last equality follows by (2.4.135).

Assume first that  $q$  is not shift invariant. Then,  $q(j_0) < \sum_k q(k, j_0)$  for some  $j_0$ . For  $u$  such that  $u(\cdot, j) = 1$  when  $j \neq j_0$  and  $u(\cdot, j_0) = e^\alpha$ ,

$$\begin{aligned} \sum_{j=1}^{|\Sigma|} \sum_{i=1}^{|\Sigma|} q(i, j) \log \left[ \frac{u(i, j)}{\left[ \sum_k u(k, i) \right] \pi(i, j)} \right] &= \sum_{j=1}^{|\Sigma|} \sum_{i=1}^{|\Sigma|} q(i, j) \log \left[ \frac{u(1, j)}{|\Sigma| u(1, i) \pi(i, j)} \right] \\ &= - \sum_{j=1}^{|\Sigma|} \sum_{i=1}^{|\Sigma|} q(i, j) \log \{ |\Sigma| \pi(i, j) \} + \alpha \left[ \sum_{i=1}^{|\Sigma|} q(i, j_0) - q(j_0) \right] \end{aligned} \quad (2.4.140)$$

which implies, by considering  $\alpha \rightarrow \infty$ , that  $I^{(2)}(q) = \infty$ .

Finally, when  $q$  is shift invariant then for any  $u \gg 0$

$$\sum_{i=1}^{|\Sigma|} \sum_{j=1}^{|\Sigma|} q(i, j) \log \left[ \frac{\sum_k u(k, i) q(j)}{\sum_k u(k, j) q(i)} \right] = 0. \quad (2.4.141)$$

Let  $u(i|j) = u(i, j) / \sum_k u(k, j)$  and  $q(i|j) = q(i, j) / q(j)$ , i.e.,  $q(i|j) = q(i, j) / \sum_k q(k, j)$  (since  $q$  is

shift invariant). Now, by (2.4.139) and (2.4.141)

$$\begin{aligned} I^{(2)}(q) - \sum_{i=1}^{|\Sigma|} q(i) H(q(\cdot|i)|\pi(i, \cdot)) &= \sup_{u \gg 0} \sum_{i=1}^{|\Sigma|} \sum_{j=1}^{|\Sigma|} q(i, j) \log \frac{u(i, j)q(i)}{[\sum_k u(k, i)]q(i, j)} \\ &= \sup_{u \gg 0} \sum_{i=1}^{|\Sigma|} \sum_{j=1}^{|\Sigma|} q(i, j) \log \frac{u(i|j)}{q(i|j)} = \sup_{u \gg 0} \left\{ - \sum_{j=1}^{|\Sigma|} q(j) H(q(\cdot|j)|u(\cdot|j)) \right\} \end{aligned} \quad (2.4.142)$$

Since  $H(q(\cdot|j)|u(\cdot|j)) \geq 0$  it follows that  $I^{(2)}(q) \leq \sum_i q(i) H(q(\cdot|i)|\pi(i, \cdot))$  with equality whenever  $q \gg 0$  (by the choice  $u = q$ ). The proof is complete for  $q$  which is not strictly positive by considering a sequence  $u_n \gg 0$  such that  $u_n \rightarrow q$  (so  $q(j)H(q(\cdot|j)|u_n(\cdot|j)) \rightarrow 0$  for each  $j$ ).  $\square$

**Exercises:**

2.4.4 Prove that for any strictly positive stochastic matrix  $\Pi$

$$J(\nu) = \inf_{\{q: \sum_i q(i, \cdot) = \nu\}} I^{(2)}(q) \quad (2.4.143)$$

where  $J(\cdot)$  is the rate function defined in (2.4.129) while  $I^{(2)}(\cdot)$  is as specified in (2.4.137).

**Hint:** There is no need to prove the above identity directly. Instead observe that the empirical measure  $L_n^X$  belongs to a set  $A$  iff  $L_n^{X, (2)} \in \{q: \sum_i q(i, \cdot) \in A\}$  (where the initial condition  $X_0 = x$  is equivalent to any initial condition  $X_{-1}X_0 = (i, x)$  for the  $\Pi^{(2)}$  chain). As the projection of any measure  $q \in M_1(\Sigma^{(2)})$  onto its marginal  $\nu \in M_1(\Sigma)$  is continuous and  $I^{(2)}(\cdot)$  controls the large deviations of  $L_n^{X, (2)}$  deduce that the right side of (2.4.143) is a rate function governing the large deviations of  $L_n^X$ . Conclude by proving the uniqueness of such a function and applying Theorem 2.4.3.

2.4.5 (a). Extend the validity of the identity (2.4.143) to any irreducible non-negative matrix  $\Pi$ .

**Hint:** First extend Theorem 2.4.4 to any irreducible stochastic matrix  $\Pi$  by replacing  $\Sigma^{(2)}$  with  $\Sigma_\Pi$  (see the remark following the statement of this theorem). Then, for any irreducible, non-negative matrix  $\Pi$  consider  $\Pi^*$  defined in exercise 2.4.3 and verify that  $I_{\Pi^*}^{(2)}(q) = I_\Pi^{(2)}(q) + \sum_i q(i) \log \phi(i)$ . (b). Deduce by applying the identities (2.4.134) and (2.4.143) that for any non-negative irreducible matrix  $\Pi$

$$-\log \rho(\Pi) = \inf_{q \in M_1(\Sigma_\Pi)} I^{(2)}(q) = \inf_{q \in M_1(\Sigma_\Pi)} \inf_{\text{shift invariant}} \sum_{i=1}^{|\Sigma|} \sum_{j=1}^{|\Sigma|} q(i, j) \log \frac{q(j|i)}{\pi(i, j)}. \quad (2.4.144)$$

This is Varadhan's characterization of the spectral radius of non-negative irreducible matrices which is extremely useful for many applications.

**2.4.6** Assume that  $X_1, \dots, X_n$  have the joint law  $P_x^\pi$  where  $\Pi$  is an irreducible stochastic matrix. Let  $T_1, \dots, T_n$  be a sequence of random variables over  $\mathcal{T} \triangleq \{1, 2, \dots, \ell\}$  which are conditionally independent given any realization  $\{X_k\}_{k=1}^n$  while  $\text{Prob}(T_k = t | X_{k-1} = i) = p(i, t)$ ,  $i \in \Sigma$ ,  $t \in \mathcal{T}$ . Construct the partial sums  $S_m = \sum_{k=1}^m T_k$  and let  $K_n$  be the stopping time where  $S_m$  first hits or exceeds the integer value  $n$ . The time weighted empirical measures  $L_n^{\mathbf{X}, \mathbf{T}}$  are defined via

$$L_n^{\mathbf{X}, \mathbf{T}}(j) = \frac{1}{n} \left[ \sum_{k=1}^{K_n-1} T_k \delta_j(X_k) + (n - S_{K_n-1}) \delta_j(X_{K_n}) \right].$$

(a). Suppose that  $p(i, \ell) > 0$  for any  $i \in \Sigma$ . Prove that  $L_n^{\mathbf{X}, \mathbf{T}}$  satisfies a large deviations principle with rate function

$$J(\nu) \triangleq \inf_{A_\nu} \left\{ \sum_{i=1}^{|\Sigma|} a(i) \sum_{j=1}^{|\Sigma|} \sum_{t=1}^{\ell} q(j, t|i) \log \frac{q(j, t|i)}{\pi(i, j)p(i, t)} \right\}$$

where

$$A_\nu \triangleq \{q(\cdot|i) \in M_1(\Sigma \times \mathcal{T}), a(i) \geq 0 : \sum_{i=1}^{|\Sigma|} a(i) \sum_{t=1}^{\ell} q(j, t|i) = a(j), \sum_{i=1}^{|\Sigma|} a(i) \sum_{t=1}^{\ell} tq(j, t|i) = \nu(j), \forall j \in \Sigma\}.$$

**Hint:** Interpret the event  $\{T_k = t\}$  as if the Markov chain freezes in its current state  $X_k$  for  $t$  time units and observe that  $L_n^{\mathbf{X}, \mathbf{T}}$  is merely the standard empirical measure in this new time scale. Consider the Markov chain whose state space  $\Sigma \times \mathcal{T}$  consists of the original states and the future time spans in which state changes are still forbidden (starting at the initial state  $(x, 1)$ ). Show that the transition from  $(i, s)$  to  $(j, t)$  in this chain has probability  $\pi(i, j)p(i, t)$  when  $s = 1$  and  $\delta_i(j)\delta_{s-1}(t)$  otherwise. Apply Theorem 2.4.4 to the pair empirical measure of this chain  $L_n^{\mathbf{X}, (2)}$ . Finally, observe that  $L_n^{\mathbf{X}, \mathbf{T}}(j) = \sum_{i,s,t} t L_n^{\mathbf{X}, (2)}((i, s), (j, t))$  for any  $j \in \Sigma$ , and apply a "contraction argument" of the type hinted about in exercise 2.4.4.

(b). Suppose that  $p(i, t) = p(t)$  and  $\pi(i, j) = \mu(j)$ . Prove that now the rate function for  $L_n^{\mathbf{X}, \mathbf{T}}$  is

$$J(\nu) = \inf_{\{q: E_q[T_1 \delta_j(X_1)] = \nu(j) E_q[T_1], \forall j\}} \frac{H(q|\mu \times p)}{E_q(T_1)}$$

where  $E_q(T_1) = \sum_{j,t} tq(j, t)$  and  $q \in M_1(\Sigma \times \mathcal{T})$ .

2.4.7 (a). Prove that

$$\frac{1}{n} \log P_x^\pi(X_1 = x_1, \dots, X_n = x_n) = \sum_{i=1}^{|\Sigma|} \sum_{j=1}^{|\Sigma|} L_n^{x, (2)}(i, j) \log \pi(i, j)$$

for any sequence  $\mathbf{x} = (x_1, \dots, x_n) \in \Sigma^n$  of non-zero  $P_x^\pi$  probability.

(b). Let

$$\mathcal{L}_n \triangleq \{q : q = L_n^{x, (2)}, P_x^\pi(X_1 = x_1, \dots, X_n = x_n) > 0 \text{ for some } \mathbf{x} \in \Sigma^n\}$$

be the set of possible types of pairs of states of the Markov chain. Prove that  $\mathcal{L}_n \subset M_1(\Sigma_\Pi)$  and  $|\mathcal{L}_n| \leq (n+1)^{|\Sigma|^2}$ .

(c). Let  $T(q)$  be the type class of  $q \in \mathcal{L}_n$ , namely the set of sequences  $\mathbf{x}$  of positive  $P_x^\pi$  probability for which  $L_n^{x, (2)} = q$  and let  $H(q) \triangleq \sum_{i,j} q(i, j) \log q(j|i)$ . Suppose that for any  $q \in \mathcal{L}_n$

$$(n+1)^{-(|\Sigma|^2+|\Sigma|)} e^{nH(q)} \leq |T(q)| \leq e^{nH(q)} \quad (2.4.145)$$

and moreover that

$$\lim_{n \rightarrow \infty} d_V(q, \mathcal{L}_n) = 0 \quad \forall q \in M_1(\Sigma_\Pi), \quad q \text{ shift invariant}. \quad (2.4.146)$$

Prove by adapting the method of types of Section 2.1.1 that  $L_n^{x, (2)}$  satisfies a large deviations principle with the rate function  $I^{(2)}(\cdot)$  specified in (2.4.137).

**Remark:** The estimates (2.4.145) and (2.4.146) are consequences of a somewhat involved combinatorial estimate of  $|T(q)|$  (see for example [18], eq. (35)-(37) and references therein)

## 2.5 Long rare segments in random walks

Let  $X_1, \dots, X_n, \dots$  be i.i.d. random vectors in  $\mathbb{R}^d$  with  $\Lambda(\cdot)$  a steep function such that  $\Lambda(\lambda) < \infty$  for some open ball around the origin. Let  $A$  be any rare  $\Lambda^*$  continuity subset of  $\mathbb{R}^d$ , namely, such that

$$I_A \triangleq \inf_{x \in A} \Lambda^*(x) = \inf_{x \in A^\circ} \Lambda^*(x) > 0, \quad (2.5.147)$$

where  $\Lambda^*(x)$  is the Legendre transform of  $\Lambda(\lambda)$  defined in (2.3.60).

Consider the random walk  $S_k = \sum_{i=1}^k X_i$ ,  $k = 1, 2, \dots$ ,  $S_0 = 0$  and let

$$R_n^{(A)} \triangleq \max \left\{ m - k : 0 \leq k < m \leq n, \frac{S_m - S_k}{m - k} \in A \right\}. \quad (2.5.148)$$

Thus,  $R_n^{(A)}$  is the maximal length among all segments of the random walk up to  $n$  in which the empirical mean is within the set  $A$ .

Associated with  $R_n^{(A)}$  is the dual variable

$$T_r^{(A)} \triangleq \inf \left\{ m : \frac{S_m - S_k}{m - k} \in A \text{ for some } 0 \leq k \leq m - r \right\}, \quad (2.5.149)$$

so that  $\{R_n^{(A)} \geq r\}$  if and only if  $\{T_r^{(A)} \leq n\}$ .

The analysis of the random variables  $R_n^{(A)}$  and  $T_r^{(A)}$  has applications in problems of the statistical analysis of DNA sequence matching and in the analysis of search algorithms in computer science. The following Theorem yields estimates on rare events which are usually associated with the probability of errors for matching algorithms. For some applications and refinements of these estimates, c.f. [23] and the exercises at the end of this section.

**Theorem 2.5.1**  $\lim_{n \rightarrow \infty} (R_n^{(A)} / \log n) = \lim_{r \rightarrow \infty} (r / \log T_r^{(A)}) = 1/I_A$ , *almost surely*.

**Proof:** By the Borel-Cantelli lemma and utilizing the duality of events  $\{R_n^{(A)} \geq r\} \equiv \{T_r^{(A)} \leq n\}$  the theorem follows from the estimates

$$\sum_{r=1}^{\infty} \text{Prob}(T_r^{(A)} \leq e^{r(I_A - \epsilon)}) < \infty, \quad \forall \epsilon > 0 \quad (2.5.150)$$

$$\sum_{r=1}^{\infty} \text{Prob}(T_r^{(A)} > e^{r(I_A + \epsilon)}) < \infty, \quad \forall \epsilon > 0 \quad (2.5.151)$$

when  $I_A < \infty$  and from

$$\sum_{r=1}^{\infty} \text{Prob}(T_r^{(A)} \leq e^{r/\epsilon}) < \infty, \quad \forall \epsilon > 0 \quad (2.5.152)$$

when  $I_A = \infty$ .

The desired estimates (2.5.150), (2.5.151), and (2.5.152), are immediate consequences of the bounds

$$\text{Prob}(T_r^{(A)} > n) \leq e^{-\lfloor \frac{n}{r} \rfloor \mu_r(A)} \quad (2.5.153)$$

and

$$\text{Prob}(T_r^{(A)} \leq n) \leq n \sum_{\ell=r}^{\infty} \mu_{\ell}(A), \quad (2.5.154)$$

coupled with Cramer's theorem (Theorem 2.3.2)

$$\lim_{\ell \rightarrow \infty} \frac{1}{\ell} \log \mu_\ell(A) = -I_A, \quad (2.5.155)$$

where  $\mu_\ell$  denotes the law of  $\hat{S}_\ell = \frac{1}{\ell} S_\ell$  for  $\ell \in \mathbb{Z}^+$  and  $\lfloor \frac{n}{r} \rfloor$  denotes the largest integer which is not larger than  $\frac{n}{r}$ . Indeed, assuming (2.5.154), one has by substituting  $n = \lfloor e^{r(I_A - \epsilon)} \rfloor$  that when  $I_A < \infty$

$$\sum_{r=1}^{\infty} \text{Prob}(T_r^{(A)} \leq e^{r(I_A - \epsilon)}) \leq \sum_{r=1}^{\infty} e^{r(I_A - \epsilon)} \sum_{\ell=r}^{\infty} c e^{-\ell(I_A - \epsilon/2)} \leq c' \sum_{r=1}^{\infty} e^{-r\epsilon/2} < \infty$$

for some positive constants  $c, c'$ . When  $I_A = \infty$  one obtains (2.5.152) by choosing  $n = \lfloor e^{r/\epsilon} \rfloor$  in (2.5.154) and following the same line of proof. Similarly, starting with (2.5.153) and choosing  $n = \lfloor e^{r(I_A + \epsilon)} \rfloor$  one obtains

$$\sum_{r=1}^{\infty} \text{Prob}(T_r^{(A)} > e^{r(I_A + \epsilon)}) \leq \sum_{r=1}^{\infty} \exp\left(-\frac{c}{r} e^{r(I_A + \epsilon)} e^{-r(I_A + \epsilon/2)}\right) \leq \sum_{r=1}^{\infty} \exp(-c'' e^{c'r}) < \infty,$$

for some positive constants  $c, c', c''$ . We turn therefore to the proof of the bounds (2.5.153) and (2.5.154). The bound (2.5.153) follows by the inclusion

$$\{T_r^{(A)} \leq n\} \supset \bigcup_{\ell=1}^{\lfloor \frac{n}{r} \rfloor} B_\ell, \quad B_\ell \triangleq \left\{ \frac{1}{r} (S_{\ell r} - S_{(\ell-1)r}) \in A \right\}, \quad (2.5.156)$$

as it implies

$$\text{Prob}(T_r^{(A)} > n) \leq 1 - \text{Prob}\left(\bigcup_{\ell=1}^{\lfloor \frac{n}{r} \rfloor} B_\ell\right) = (1 - \text{Prob}(B_1))^{\lfloor \frac{n}{r} \rfloor} \leq e^{-\lfloor \frac{n}{r} \rfloor \text{Prob}(B_1)} \quad (2.5.157)$$

due to the fact that  $\{B_\ell\}_{\ell=1}^{\infty}$  are independent events related to disjoint segments of the random walk, of equal probabilities  $\text{Prob}(B_\ell) = \mu_r(A)$ ,  $\ell = 1, 2, \dots$ .

The bound (2.5.154) follows by the inclusion

$$\{T_r^{(A)} \leq n\} \subset \bigcup_{k=0}^{n-r} \bigcup_{m=k+r}^n C_{k,m} \subset \bigcup_{k=0}^{n-1} \bigcup_{m=k+r}^{\infty} C_{k,m}, \quad C_{k,m} \triangleq \left\{ \frac{S_m - S_k}{m - k} \in A \right\}, \quad (2.5.158)$$

and the union of events bound. Note that  $\text{Prob}(C_{k,m}) = \mu_{m-k}(A)$ , and  $m - k \geq r$  while in (2.5.158), there are at most  $n$  possible choices of  $k$ .  $\square$

**Remark:** Note that Theorem 2.5.1 holds as long as (2.5.155) holds. For example, consider exercises 2.5.3 and 2.5.4.

**Exercises:**



2.5.1 Suppose that  $I_A < \infty$  and (2.5.155) may be refined to

$$\lim_{r \rightarrow \infty} [\mu_r(A) r^{d/2} e^{r I_A}] = a ,$$

for some  $a \in (0, \infty)$  (such an example is presented in Section 2.10 for  $d = 1$ ). Let

$$\hat{R}_n^{(A)} \triangleq \frac{I_A R_n^{(A)} - \log n}{\log \log n} + \frac{d}{2} .$$

- (a). Repeat the above calculations and prove that  $\limsup_{n \rightarrow \infty} |\hat{R}_n^{(A)}| \leq 1$  almost surely.  
(b). Also deduce that  $\lim_{n \rightarrow \infty} \text{Prob}(\hat{R}_n^{(A)} \leq -\epsilon) = 0$  for all  $\epsilon > 0$ .

2.5.2 Let  $A = \{1\}$  and  $X_i$  be i.i.d. Bernoulli( $p$ ) random variables. Then,  $R_n^{(A)}$  is the longest consecutive run of 1-s in the binary sequence  $X_1, \dots, X_n$ . Let the renewal times  $Z_1, Z_2, \dots$  be the locations of zeros in this sequence (with  $Z_0 \triangleq 0$ ). Then,  $Q_k \triangleq Z_k - Z_{k-1} - 1$ ,  $k = 1, 2, \dots$  are i.i.d. Geometric( $1 - p$ ) random variables and  $R_{Z_k}^{(A)} = \max\{Q_1, \dots, Q_k\}$ . By standard renewal theory  $\frac{Z_k}{k} \rightarrow \frac{1}{1-p}$  as  $k \rightarrow \infty$ , almost surely. Verify that here  $I_A = -\log p$  and deduce that

$$\lim_{n \rightarrow \infty} \epsilon_n (I_A R_n^{(A)} - \log n) = 0$$

almost surely, whenever  $\lim_{n \rightarrow \infty} \epsilon_n = 0$ .

2.5.3 (a). Consider a sequence  $X_1, \dots, X_n$  of i.i.d. random variables over a finite alphabet  $\Sigma$  having a marginal law  $\mu$  such that  $\Sigma_\mu = \Sigma$ . Let  $R_n^{(\Gamma)}$  be the longest among all segments of this sequence with segmental empirical measures in the open set  $\Gamma \subset M_1(\Sigma)$ . Assume that  $\mu \notin \bar{\Gamma}$  and derive the analog of Theorem 2.5.1 for this situation.

(b). Assume further that  $\Gamma$  is convex. Let  $\nu^*$  be the unique minimizer of  $H(\cdot | \mu)$  in  $\bar{\Gamma}$ . Prove that as  $n \rightarrow \infty$  the empirical measures associated with the segments contributing to  $R_n^{(\Gamma)}$  converge almost surely to  $\nu^*$ .

Hint: Let  $\Gamma_\delta \triangleq \Gamma \cap \overline{B_{\nu^*, \delta}^c}$  where  $B_{\nu^*, \delta}$  is an open ball of radius  $\delta > 0$  around  $\nu^*$ . Prove that  $\limsup_{n \rightarrow \infty} \frac{R_n^{(\Gamma_\delta)}}{R_n^{(\Gamma)}} < 1$  almost surely, for any  $\delta > 0$ . Deduce that for any  $\delta > 0$  the empirical measures of the segments contributing to  $R_n^{(\Gamma)}$  are eventually within distance  $\delta$  of  $\nu^*$  almost surely.

2.5.4 Assume that  $X_1, \dots, X_n$  and  $Y_1, \dots, Y_n$  are as in exercise 2.4.1. Specifically,  $X_k$  are the states of a Markov chain over the finite set  $\{1, 2, \dots, |\Sigma|\}$  with an irreducible transition matrix  $\Pi$ ,

and the conditional law of  $Y_k$  when  $X_k = j$  is  $\mu_j \in M_1(\mathbb{R}^d)$  while  $\{Y_k\}$  are independent given any realization of the Markov chain states. Further, suppose that the logarithmic moment generating functions  $\Lambda_j$  associated with  $\mu_j$  are finite everywhere and define the matrices  $\Pi_\lambda$  via

$$\pi_\lambda(i, j) \triangleq \pi(i, j) e^{\Lambda_j(\lambda)}.$$

Let  $\Lambda^*(x)$  denote the Legendre transform of  $\log \rho(\Pi_\lambda)$  and suppose  $A$  is a rare  $\Lambda^*$  continuity set. Let  $R_n^{(A)}$  be the longest among all segments whose  $Y$ -segmental empirical mean belongs to  $A \subset \mathbb{R}^d$ . Prove that Theorem 2.5.1 holds with  $\Lambda^*$  as defined here.

## 2.6 The Gibbs conditioning principle in finite alphabet

Let  $X_1, X_2, \dots, X_n$  be a sequence of i.i.d. random variables with law  $\mu$  over the finite alphabet  $\Sigma \subset \mathbb{R}^1$  and assume without loss of generality that  $\Sigma_\mu = \Sigma$ . The following question is of fundamental importance in statistical mechanics. Given a set  $A \in \mathbb{R}$  and a constraint of the type  $\hat{S}_n \in A$  what is the conditional law of  $X_1$  for large  $n$ ? In other words, what are the limit points of the conditional probability vector

$$\mu_n^*(a_i) \triangleq \text{Prob}_\mu(X_1 = a_i | \hat{S}_n \in A) \quad i = 1, \dots, |\Sigma|, \quad (2.6.159)$$

as  $n \rightarrow \infty$  (recall that  $\hat{S}_n \triangleq \frac{1}{n} \sum_{j=1}^n X_j = \langle L_n^X, \mathbf{a} \rangle$ ). Note that for any function  $f : \Sigma \rightarrow \mathbb{R}^1$

$$\begin{aligned} \langle \mu_n^*, f \rangle &= E[f(X_1) | \hat{S}_n \in A] = E[f(X_2) | \hat{S}_n \in A] \\ &= E\left[\frac{1}{n} \sum_{j=1}^n f(X_j) | \hat{S}_n \in A\right] = E[\langle L_n^X, f \rangle | \langle L_n^X, \mathbf{a} \rangle \in A] \end{aligned} \quad (2.6.160)$$

where we have used the fact that  $X_j$  are identically distributed (although not independent) even under the conditioning  $\hat{S}_n \in A$ . Therefore,

$$\mu_n^* = E[L_n^X | L_n^X \in \Gamma], \quad (2.6.161)$$

where  $\Gamma \triangleq \{\nu : \langle \nu, \mathbf{a} \rangle \in A\}$  (compare with (2.1.19) in Section 2.1.2). Therefore, the characterization of possible limit points of the sequence  $\mu_n^*$  as  $n \rightarrow \infty$  can be cast in terms of conditional limits for the empirical measures  $L_n^X$ .

The following characterization of the limits of  $\mu_n^*$  is a consequence of Theorem 2.1.1 for any non-empty set  $\Gamma$  which is an  $H(\cdot|\mu)$  continuity set, namely,

$$I_\Gamma \triangleq \inf_{\nu \in \Gamma^\circ} H(\nu|\mu) = \inf_{\nu \in \bar{\Gamma}} H(\nu|\mu) . \quad (2.6.162)$$

**Theorem 2.6.1 (Gibb's principle)**

(a). *The set of possible limit points of  $\mu_n^*$  is the closure of the convex hull of*

$$\mathcal{M}_\Gamma \triangleq \{\nu \in \bar{\Gamma} : H(\nu|\mu) = I_\Gamma\} \quad (2.6.163)$$

(b). *For any convex set  $\Gamma$  of non-empty interior the set  $\mathcal{M}_\Gamma$  is a point to which  $\mu_n^*$  converges as  $n \rightarrow \infty$ .*

**Remark:** For conditions on  $\Gamma$  (alternatively, on  $A$ ) under which (2.6.162) holds see exercises 2.1.1-2.1.3 in Section 2.1.1.

**Proof:** As  $|\Sigma| < \infty$ ,  $\bar{\Gamma}$  is a compact set and thus  $\mathcal{M}_\Gamma$  is non-empty. Moreover, part (b) of the theorem follows from part (a) by exercise 2.1.3 and the compactness of  $M_1(\Sigma)$  (in that exercise you showed that indeed (2.6.162) holds when  $\Gamma$  is a convex set of non-empty interior and that the set  $\mathcal{M}_\Gamma$  is a point).

We shall prove that for any  $\delta > 0$

$$\lim_{n \rightarrow \infty} \text{Prob}(L_n^{\mathbf{X}} \in \mathcal{M}_\Gamma^\delta | L_n^{\mathbf{X}} \in \Gamma) = 1 , \quad (2.6.164)$$

with an exponential (in  $n$ ) rate of convergence, where  $\mathcal{M}_\Gamma^\delta \triangleq \{\nu : d_V(\nu, \mathcal{M}_\Gamma) < \delta\}$ .

Since  $M_1(\Sigma)$  is a bounded set, (2.6.161) and (2.6.164) imply that for any  $\delta > 0$ ,  $\mu_n^*$  eventually belongs to the convex hull of  $\mathcal{M}_\Gamma^{2\delta}$ . All points in the convex hull of  $\mathcal{M}_\Gamma^\delta$  are within variational distance  $\delta$  of some point in the convex hull of  $\mathcal{M}_\Gamma$  (since  $d_V$  is a convex function on  $M_1(\Sigma) \times M_1(\Sigma)$ ). Thus, since  $\delta$  is arbitrarily small, limit points of  $\mu_n^*$  are necessarily in the closure of the convex hull of  $\mathcal{M}_\Gamma$  as claimed.

The limit (2.6.164) definitely follows from

$$\begin{aligned} \limsup_{n \rightarrow \infty} \frac{1}{n} \log \text{Prob}_\mu(L_n^{\mathbf{X}} \in (\mathcal{M}_\Gamma^\delta)^c | L_n^{\mathbf{X}} \in \Gamma) &= \\ \limsup_{n \rightarrow \infty} \left\{ \frac{1}{n} \log \text{Prob}_\mu(L_n^{\mathbf{X}} \in (\mathcal{M}_\Gamma^\delta)^c \cap \Gamma) - \frac{1}{n} \log \text{Prob}_\mu(L_n^{\mathbf{X}} \in \Gamma) \right\} &< 0 . \end{aligned} \quad (2.6.165)$$

However, by Theorem 2.1.1 and (2.6.162)

$$I_\Gamma = - \lim_{n \rightarrow \infty} \frac{1}{n} \log \text{Prob}_\mu(L_n^{\mathbf{X}} \in \Gamma) , \quad (2.6.166)$$

whereas by Theorem 2.1.1 also

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log \text{Prob}_\mu(L_n^{\mathbf{X}} \in (\mathcal{M}_\Gamma^\delta)^c \cap \Gamma) \leq - \inf_{\nu \in (\mathcal{M}_\Gamma^\delta)^c \cap \bar{\Gamma}} H(\nu|\mu) . \quad (2.6.167)$$

Observe that  $\mathcal{M}_\Gamma^\delta$  are open sets and therefore  $(\mathcal{M}_\Gamma^\delta)^c \cap \bar{\Gamma}$  are compact sets. Thus, for some  $\tilde{\nu} \in (\mathcal{M}_\Gamma^\delta)^c \cap \bar{\Gamma}$

$$\inf_{\nu \in (\mathcal{M}_\Gamma^\delta)^c \cap \bar{\Gamma}} H(\nu|\mu) = H(\tilde{\nu}|\mu) > I_\Gamma , \quad (2.6.168)$$

where the above inequality follows from the definition of  $\mathcal{M}_\Gamma$  since  $\tilde{\nu} \notin \mathcal{M}_\Gamma$  while  $\tilde{\nu} \in \bar{\Gamma}$ . Finally, (2.6.164) follows from (2.6.165)-(2.6.168).  $\square$

### Remarks

- (a). Intuitively one expects  $X_1, \dots, X_k$  to be asymptotically independent (as  $n \rightarrow \infty$ ) for any fixed  $k$ , when the conditioning event is  $\{L_n^{\mathbf{X}} \in \Gamma\}$ . This is indeed shown in exercise 2.6.3 by considering “super-symbols” from the enlarged alphabet  $\Sigma^k$ .
- (b). Theorem 2.6.1 holds for any set  $\Gamma$  satisfying (2.6.162). However, the particular conditioning set  $\{\nu : \langle \nu, \mathbf{a} \rangle \in A\}$  has an important significance in statistical mechanics because it represents an energy-like constraint.
- (c). Recall the relationship (2.1.23) of Section 2.1.2 which implies that for any non-empty, convex, open set  $A \subset K$  the unique limit of  $\mu_n^*$  is of the form

$$\nu_\lambda(a_i) = e^{\lambda a_i - \Lambda(\lambda)} \mu(a_i) ,$$

for some appropriately chosen  $\lambda \in \mathbb{R}^1$  which is called the **Gibbs parameter** associated with  $A$ . In particular, for any  $x \in K^\circ$  the Gibbs parameter  $\lambda$  associated with  $A = (x - \delta, x + \delta)$  converges as  $\delta \rightarrow 0$  to  $\lambda^{(x)}$ , the unique solution of the equation  $\Lambda'(\lambda) = x$  (for details see Section 2.1.2).

- (d). A Gibbs conditioning principle holds beyond the i.i.d. case. Actually, all that is needed is that  $X_i$  are exchangeable conditionally upon any given value of  $L_n^{\mathbf{X}}$  (so that (2.6.161) holds). For such an example, consider exercise 2.6.2.

Exercises:

2.6.1 Prove Theorem 2.6.1 by the method of types (specifically, use Lemma 2.1.4 directly).

2.6.2 Prove the Gibbs conditioning principle for sampling without replacement.

- (a). Observe that again  $X_j$  are identically distributed even when  $L_n^X$  is given. Conclude that (2.6.161) holds.
- (b). Assume that  $\Gamma$  is such that

$$I_\Gamma = \inf_{\nu \in \Gamma^o} I_{\beta, \mu}(\nu) = \inf_{\nu \in \bar{\Gamma}} I_{\beta, \mu}(\nu) < \infty ,$$

and define  $\mathcal{M}_\Gamma \triangleq \{\nu \in \bar{\Gamma} : I_{\beta, \mu}(\nu) = I_\Gamma\}$ . Prove that now both parts of Theorem 2.6.1 hold (for part (b) you may rely on exercise 2.1.12).

2.6.3 (a). Suppose that  $\Sigma = (\Sigma')^k$  and  $\mu = (\mu')^k$  are a  $k$ -th product alphabet and a  $k$ -th product underlying measure on it and assume that  $\Sigma'_{\mu'} = \Sigma'$  (as usual). For any law  $\nu \in M_1(\Sigma)$  let  $\nu^{(j)} \in M_1(\Sigma')$ ,  $j = 1, \dots, k$  denote its  $j$ -th marginal on  $\Sigma'$ . Prove that

$$\frac{1}{k} H(\nu | \mu) \geq \frac{1}{k} \sum_{j=1}^k H(\nu^{(j)} | \mu') \geq H\left(\frac{1}{k} \sum_{j=1}^k \nu^{(j)} | \mu'\right) ,$$

with equality if and only if  $\nu = (\nu')^k$  for some  $\nu' \in M_1(\Sigma')$ .

(b). Assume that

$$\Gamma \triangleq \left\{ \nu : \frac{1}{k} \sum_{j=1}^k \nu^{(j)} \in \Gamma' \right\} \quad (2.6.169)$$

for some  $\Gamma' \subset M_1(\Sigma')$  which satisfies (2.6.162) with respect to  $\mu'$ . Prove that then  $\mathcal{M}_\Gamma = (\mathcal{M}_{\Gamma'})^k$  and conclude that any limit point of  $\mu_n^*$  is a  $k$ -th product of some appropriate law on  $\Sigma'$ .

(c). Consider now the  $k$ -th joint conditional law

$$\mu_n^*(a'_{i_1}, \dots, a'_{i_k}) \triangleq \text{Prob}_{\mu'}(X_1 = a'_{i_1}, \dots, X_k = a'_{i_k} | L_n^X \in \Gamma') \quad a'_{i_j} \in \Sigma', j = 1, \dots, k ,$$

where  $X_i$  are i.i.d. with marginal law  $\mu' \in M_1(\Sigma')$  over the finite alphabet  $\Sigma'$  and  $\Gamma' \subset M_1(\Sigma')$  satisfies (2.6.162). Let  $\mu = (\mu')^k$  be the law of  $Y_i = (X_{1+k(i-1)}, \dots, X_{ki})$  over a new alphabet  $\Sigma$ . Prove that for any  $m \in \mathbb{Z}^+$

$$\mu_{mk}^*(a_i) = \text{Prob}_\mu(Y_1 = a_i | L_m^Y \in \Gamma) , \quad \forall a_i \in \Sigma ,$$

where  $\Gamma$  is defined in (2.6.169). Deduce that as  $n \rightarrow \infty$  along integer multiples of  $k$  the random variables  $X_i$ ,  $i = 1, \dots, k$  are asymptotically conditionally i.i.d (i.e., any limit point of  $\mu_n^*$  is a  $k$ -th product of an element of  $M_1(\Sigma')$ ).

(d). Prove that the above conclusion extends to  $n$  which need not be an integer multiple of  $k$  whenever  $\mathcal{M}_{\Gamma'}$  is a single point.

## 2.7 The hypothesis test problem

Consider the problem of hypothesis testing between two product measures for the i.i.d. random variables  $Y_1, \dots, Y_n, \dots$ . Specifically,  $Y_j$  are either distributed according to the law  $\mu_0 \in M_1(\Sigma)$  (hypothesis  $H_0$ ) or according to  $\mu_1 \in M_1(\Sigma)$  (hypothesis  $H_1$ ). The alphabet  $\Sigma$  may in general be quite arbitrary provided that the probability measures  $\mu_0$  and  $\mu_1$  are well defined (Markov chains over finite alphabet are considered in exercise 2.7.4).

**Definition 2.7.1** A decision test  $\mathcal{S}$  is a sequence of maps  $\mathcal{S}^n : \Sigma^n \rightarrow \{0, 1\}$ , for  $n = 1, 2, \dots$ , with the interpretation that when  $Y_1 = y_1, \dots, Y_n = y_n$  is observed then  $H_0$  is accepted ( $H_1$  rejected) if  $\mathcal{S}^n(y_1, \dots, y_n) = 0$  while  $H_1$  is accepted ( $H_0$  rejected) if  $\mathcal{S}^n(y_1, \dots, y_n) = 1$ .

The performance of a decision test  $\mathcal{S}$  is determined by the error probabilities

$$\alpha_n \triangleq \text{Prob}_{\mu_0}(H_0 \text{ rejected by } \mathcal{S}^n), \quad \beta_n \triangleq \text{Prob}_{\mu_1}(H_1 \text{ rejected by } \mathcal{S}^n), \quad n \in \mathbb{Z}^+. \quad (2.7.170)$$

One wishes to minimize  $\beta_n$ . If no constraint is put on  $\alpha_n$ , then one may have  $\beta_n = 0$  with the test  $\mathcal{S}^n(y_1, \dots, y_n) \equiv 1$  at the cost of  $\alpha_n = 1$ . Thus, a sensible criterion for optimality is to seek a test which minimizes  $\beta_n$  subject to a constraint on  $\alpha_n$ . Suppose now that the probability measures  $\mu_0, \mu_1$  are known a-priori and that they are *equivalent measures*, so the likelihood ratios  $L_{0||1}(y) = \frac{d\mu_0}{d\mu_1}(y)$  and  $L_{1||0}(y) = \frac{d\mu_1}{d\mu_0}(y)$  exist (some extensions for  $\mu_0, \mu_1$  which are not equivalent are given in exercise 2.7.3). This equivalence assumption is valid for example when  $\mu_0, \mu_1$  are discrete measures with  $\Sigma_{\mu_0} = \Sigma_{\mu_1}$ , or when  $\Sigma = \mathbb{R}^d$  and both  $\mu_0$  and  $\mu_1$  possess strictly positive densities. In order to avoid trivialities it is further assumed that  $\mu_0$  and  $\mu_1$  are distinguishable, i.e. that they differ on a set whose probability under  $\mu_0$  (and  $\mu_1$ ) is positive.

Let  $X_j \triangleq \log L_{1||0}(Y_j) = -\log L_{0||1}(Y_j)$  denote the observed log-likelihood ratios. These are bona-fide i.i.d. random variables over  $\mathbb{R}^1$  which are non-zero with positive probability. Moreover,

$$\bar{x}_0 \triangleq E_{\mu_0}[X_1] = E_{\mu_1}[X_1 e^{-X_1}] ,$$

exists (with possibly  $\bar{x}_0 = -\infty$ ) as  $xe^{-x} \leq 1$ . Similarly,

$$\bar{x}_1 \triangleq E_{\mu_1}[X_1] = E_{\mu_0}[X_1 e^{X_1}] > E_{\mu_0}[X_1] = \bar{x}_0 ,$$

exists (with possibly  $\bar{x}_1 = \infty$ ) and the above inequality is strict since  $X_1$  is non-zero with positive probability. Note that  $\bar{x}_0$  and  $\bar{x}_1$  may be both characterized in terms of relative entropy, c.f. exercise 2.7.1.

**Definition 2.7.2** A Neyman-Pearson test is a test in which for any  $n$  the mean observed log-likelihood ratio  $\hat{S}_n \triangleq \frac{1}{n} \sum_{j=1}^n X_j$  is compared against a threshold  $\gamma_n$  and  $H_1$  is accepted (rejected) when  $\hat{S}_n > \gamma_n$  ( $\hat{S}_n \leq \gamma_n$ ).

It is well known that Neyman-Pearson tests are optimal in the sense that there are neither tests with the same value of  $\alpha_n$  and a smaller value of  $\beta_n$  nor tests with the same value of  $\beta_n$  and a smaller value of  $\alpha_n$  (see for example [7] for a simple proof of this claim).

The exponential rates of  $\alpha_n$  and  $\beta_n$  for Neyman-Pearson tests with constant thresholds  $\gamma \in (\bar{x}_0, \bar{x}_1)$  are thus of particular interest. These may be cast in terms of the large deviations of  $\hat{S}_n$ . In particular, since  $X_j$  are i.i.d. real valued random variables, the following theorem is an application of Theorem 2.2.1.

**Theorem 2.7.1** For any Neyman-Pearson test with constant threshold  $\gamma \in (\bar{x}_0, \bar{x}_1)$

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \alpha_n = -\Lambda_0^*(\gamma) < 0 , \quad (2.7.171)$$

and

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \beta_n = \gamma - \Lambda_0^*(\gamma) < 0 , \quad (2.7.172)$$

where

$$\Lambda_0^*(x) \triangleq \sup_{\lambda \in [0,1]} \{ \lambda x - \Lambda_0(\lambda) \} \quad (2.7.173)$$

and  $\Lambda_0(\lambda) = \log E_{\mu_0}[e^{\lambda X_1}]$  is the logarithmic moment generating function of  $X_1$  under  $H_0$ .

**Proof:** Both (2.7.171) and (2.7.172) follow by slight modification of the proof of Theorem 2.2.1. First note that

$$\bar{x}_0 = \lim_{\lambda \downarrow 0} \Lambda'_0(\lambda) < \gamma < \lim_{\lambda \uparrow 1} \Lambda'_0(\lambda) = \bar{x}_1.$$

Thus,  $\gamma = \Lambda'_0(\eta)$  for some  $\eta \in (0, 1)$  and  $\Lambda_0^*(\gamma)$  indeed equals the Legendre transform of  $\Lambda_0$  at the point  $\gamma$ . Now by (2.2.46)

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log \alpha_n = \limsup_{n \rightarrow \infty} \frac{1}{n} \log \text{Prob}_{\mu_0}(\hat{S}_n \in (\gamma, \infty)) \leq - \inf_{x \geq \gamma} \Lambda_0^*(x) = -\Lambda_0^*(\gamma) \quad (2.7.174)$$

where the last equality follows since  $\Lambda_0^*(\cdot)$  is nondecreasing on  $[\gamma, \infty)$  as  $\gamma > \bar{x}_0$ .

By the definition of  $X_j$  the logarithmic moment generating function associated with  $\mu_1$  is merely  $\Lambda_0(\lambda + 1)$  and so the rate function  $\Lambda_1^*(x) \triangleq \Lambda_0^*(x) - x$  governs the large deviations bounds for the laws of  $\hat{S}_n$  under  $H_1$ . Apply (2.2.46) once again to obtain

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log \beta_n = \limsup_{n \rightarrow \infty} \frac{1}{n} \log \text{Prob}_{\mu_1}(\hat{S}_n \in (-\infty, \gamma]) \leq - \inf_{x \leq \gamma} \Lambda_1^*(x) = -\Lambda_1^*(\gamma) \quad (2.7.175)$$

where the last equality follows since  $\Lambda_1^*(\cdot)$  is nonincreasing on  $(-\infty, \gamma]$  as  $\gamma < \bar{x}_1$ .

Since  $\gamma = \Lambda'_0(\eta)$  for some  $\eta \in (0, 1)$  where  $\Lambda_0$  is a strictly convex,  $C^\infty$  function, both  $\Lambda_0^*(\cdot)$  and  $\Lambda_1^*(\cdot)$  are continuous at the point  $\gamma$  (consider further exercise 2.2.5 of Section 2.2). Moreover, for large enough  $r$  the lower bound (2.2.51) applies to  $y_r = \gamma + \frac{1}{r}$  and  $\delta_r = \frac{1}{r}$  implying that

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log \alpha_n \geq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \text{Prob}_{\mu_0}(\hat{S}_n \in B_{y_r, \delta_r}) \geq -\Lambda_0^*(y_r) \quad (2.7.176)$$

where  $B_{y, \delta} \triangleq (y - \delta, y + \delta)$ . By taking the limit  $r \rightarrow \infty$  and combining the above lower bound with (2.7.174) one deduces (2.7.171). Similarly, one has for  $z_r = \gamma - \frac{1}{r}$  (with  $\delta_r = \frac{1}{r}$  and  $r$  large)

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log \beta_n \geq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \text{Prob}_{\mu_1}(\hat{S}_n \in B_{z_r, \delta_r}) \geq -\Lambda_1^*(z_r) \quad (2.7.177)$$

implying (2.7.172) in the limit  $r \rightarrow \infty$ . □

**Remarks:** (a). Observe that Theorem 2.7.1 holds even when  $\bar{x}_0 = -\infty$  or  $\bar{x}_1 = \infty$  or both. Its proof is actually a specialization of exercise 2.2.7 from Section 2.2.

(b). A refinement of Theorem 2.7.1 is given in exercise 2.10.3 where the exact limiting behavior of  $\alpha_n$  ( $\beta_n$ ) is derived.

A corollary of Theorem 2.7.1 is Chernoff's asymptotic bound on the best achievable Bayesian probability of error

$$P_n^{(e)} = \text{Prob}(H_0)\alpha_n + \text{Prob}(H_1)\beta_n \quad (2.7.178)$$



**Corollary 2.7.1 (Chernoff's Bound)** *If  $0 < \text{Prob}(H_0) < 1$  then*

$$\liminf_{n \rightarrow \infty} \left\{ \frac{1}{n} \log P_n^{(e)} \right\} = -\Lambda_0^*(0) = - \inf_{\lambda \in [0,1]} \Lambda_0(\lambda) \quad , \quad (2.7.179)$$

*where the above infimum is over all tests.*

**Remarks:**

(a). In particular, Theorem 2.7.1 thus implies that the best Bayesian exponential error rate is achieved by a Neyman-Pearson test with *zero threshold*.

(b). The rate  $\Lambda_0^*(0)$  is called Chernoff's information of the measures  $\mu_0$  and  $\mu_1$ .

**Proof:** It suffices to consider only Neyman-Pearson tests. Let  $\alpha_n^*$  and  $\beta_n^*$  be the error probabilities for the zero threshold Neyman-Pearson test. Then by (2.7.171) and (2.7.172)

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \alpha_n^* = \lim_{n \rightarrow \infty} \frac{1}{n} \log \beta_n^* = -\Lambda_0^*(0) \quad (2.7.180)$$

For any test either  $\alpha_n \geq \alpha_n^*$  (when  $\gamma_n \leq 0$ ), or  $\beta_n \geq \beta_n^*$  (when  $\gamma_n \geq 0$ ). Thus, for any test

$$\frac{1}{n} \log P_n^{(e)} \geq \frac{1}{n} \log [\min\{\text{Prob}(H_0), \text{Prob}(H_1)\}] + \min\left\{ \frac{1}{n} \log \alpha_n^*, \frac{1}{n} \log \beta_n^* \right\} .$$

As  $0 < \text{Prob}(H_0) < 1$ , the limit  $n \rightarrow \infty$  yields (2.7.179) in view of (2.7.180). □

Another corollary of Theorem 2.7.1 is the following lemma which determines the best exponential rate for  $\beta_n$  when  $\alpha_n$  are bounded away from 1.

**Lemma 2.7.1 (Stein's Lemma)** *Let  $\beta_n^\epsilon$  be the minimum of  $\beta_n$  among all tests with  $\alpha_n < \epsilon$ . Then, for any  $\epsilon < 1$*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \beta_n^\epsilon = \bar{x}_0 \quad . \quad (2.7.181)$$

**Proof:** It clearly suffices to consider only Neyman-Pearson tests. Then,

$$\alpha_n = \text{Prob}_{\mu_0}(\hat{S}_n > \gamma_n) \quad , \quad (2.7.182)$$

and

$$\beta_n = \text{Prob}_{\mu_1}(\hat{S}_n \leq \gamma_n) = E_{\mu_1}[1_{\hat{S}_n \leq \gamma_n}] = E_{\mu_0}[1_{\hat{S}_n \leq \gamma_n} e^{n\hat{S}_n}] , \quad (2.7.183)$$

where the last equality follows from the definition of  $X_j$  (as the observed log-likelihood ratios). The identity (2.7.183) yields

$$\frac{1}{n} \log \beta_n = \frac{1}{n} \log E_{\mu_0}[1_{\hat{S}_n \leq \gamma_n} e^{n\hat{S}_n}] \leq \gamma_n \quad (2.7.184)$$

(a). Suppose first that  $\bar{x}_0 = -\infty$ . Then, for any Neyman-Pearson test with a fixed threshold  $\gamma$ , eventually  $\alpha_n < \epsilon$  by (2.7.171). Thus,  $\frac{1}{n} \log \beta_n^\epsilon \leq \gamma$  for any  $\gamma$  and  $n$  large enough by (2.7.184) and (2.7.181) follows.

(b). Assume now that  $\bar{x}_0 > -\infty$ . Similarly, apply (2.7.171) to deduce that eventually  $\alpha_n < \epsilon$  for Neyman-Pearson tests with a constant threshold  $\gamma > \bar{x}_0$  and so by (2.7.184)

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log \beta_n^\epsilon \leq \bar{x}_0 + \eta, \quad (2.7.185)$$

for any  $\eta > 0$  and any  $\epsilon > 0$ .

Moreover, without loss of generality one may assume that

$$\liminf_{n \rightarrow \infty} \gamma_n \geq \bar{x}_0, \quad (2.7.186)$$

for otherwise, by the weak law of large numbers  $\limsup_{n \rightarrow \infty} \alpha_n = 1$ . When (2.7.186) holds and  $\alpha_n < \epsilon$  then by (2.7.182) and the weak law of large numbers

$$\liminf_{n \rightarrow \infty} \text{Prob}_{\mu_0}(\hat{S}_n \in [\bar{x}_0 - \eta, \gamma_n]) \geq 1 - \epsilon \quad \text{for any } \eta > 0. \quad (2.7.187)$$

Hence, by (2.7.183)

$$\begin{aligned} \frac{1}{n} \log \beta_n &\geq \frac{1}{n} \log E_{\mu_0}[1_{\hat{S}_n \in [\bar{x}_0 - \eta, \gamma_n]} e^{n\hat{S}_n}] \\ &\geq \bar{x}_0 - \eta + \frac{1}{n} \log \text{Prob}_{\mu_0}(\hat{S}_n \in [\bar{x}_0 - \eta, \gamma_n]). \end{aligned} \quad (2.7.188)$$

By combining (2.7.187), (2.7.188) and the optimality of these Neyman-Pearson tests one obtains

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log \beta_n^\epsilon \geq \bar{x}_0 - \eta \quad \text{for any } \eta > 0. \quad (2.7.189)$$

The desired limit (2.7.181) is now a direct consequence of (2.7.185) and (2.7.189).  $\square$

### Exercises:

**2.7.1** Suppose that  $Y_1, \dots, Y_n$  are i.i.d. random variables over the finite set  $\Sigma \triangleq \{a_1, \dots, a_{|\Sigma|}\}$  and  $\Sigma_{\mu_0} = \Sigma_{\mu_1} = \Sigma$  (namely,  $\mu_0$  and  $\mu_1$  are strictly positive over  $\Sigma$ ).

(a). Prove that  $\bar{x}_1 = H(\mu_1|\mu_0) < \infty$  and  $\bar{x}_0 = -H(\mu_0|\mu_1) > -\infty$  (see Section 2.1 for the definition

of relative entropy).

(b). For  $\eta \in [0, 1]$  define the probability measures

$$\mu_\eta(a_j) \triangleq \frac{\mu_1(a_j)^\eta \mu_0(a_j)^{1-\eta}}{\sum_{k=1}^{|\Sigma|} \mu_1(a_k)^\eta \mu_0(a_k)^{1-\eta}} \quad j = 1, \dots, |\Sigma|.$$

Let  $\gamma = H(\mu_\eta|\mu_0) - H(\mu_\eta|\mu_1)$  and prove that  $\Lambda_0^*(\gamma) = H(\mu_\eta|\mu_0)$ .

**2.7.2** Consider the scenario of exercise 2.7.1.

(a). Define the conditional probability vectors

$$\mu_n^*(a_j) \triangleq \text{Prob}_{\mu_0}(Y_1 = a_j \mid H_0 \text{ rejected by } S^n) \quad j = 1, \dots, |\Sigma|, \quad (2.7.190)$$

where  $S$  is a Neyman-Pearson test with a fixed threshold  $\gamma \triangleq H(\mu_\eta|\mu_0) - H(\mu_\eta|\mu_1)$  and  $\eta \in (0, 1)$ .

Use Theorem 2.6.1 to deduce that  $\mu_n^* \rightarrow \mu_\eta$  as  $n \rightarrow \infty$ .

Hint: You may also find parts of the proof of Theorem 2.1.2 useful for solving this problem.

(b). Consider now the  $k$ -th joint conditional law

$$\mu_n^*(a_{j_1}, \dots, a_{j_k}) \triangleq \text{Prob}_{\mu_0}(Y_1 = a_{j_1}, \dots, Y_k = a_{j_k} \mid H_0 \text{ rejected by } S^n) \quad a_{j_\ell} \in \Sigma, \ell = 1, \dots, k.$$

Apply exercise 2.6.3 in order to deduce that

$$\lim_{n \rightarrow \infty} \mu_n^*(a_{j_1}, \dots, a_{j_k}) = \mu_\eta(a_{j_1}) \mu_\eta(a_{j_2}) \cdots \mu_\eta(a_{j_k}) \quad .$$

Try to interpret this result.

**2.7.3** Suppose that  $L_{1||0}(y) = \frac{d\mu_1}{d\mu_0}(y)$  does not exist while  $L_{0||1}(y) = \frac{d\mu_0}{d\mu_1}(y)$  does exist. Prove that Stein's lemma holds true whenever  $\bar{x}_0 \triangleq -E_{\mu_0}[\log L_{0||1}(Y_1)] > -\infty$ .

Hint: Split  $\mu_1$  into its singular part with respect to  $\mu_0$  and its restriction on the support of the measure  $\mu_0$ .

**2.7.4** Suppose that  $Y_1, \dots, Y_n$  are the states of a Markov chain over the finite set  $\Sigma = \{1, 2, \dots, |\Sigma|\}$  where the initial state of the chain  $Y_0$  is known a-priori to be some  $x \in \Sigma$ . The transition matrix under  $H_0$  is  $\Pi_0$  while under  $H_1$  it is  $\Pi_1$ , both of which are irreducible matrices with the same set of non-zero values. Here the Neyman-Pearson tests are based upon  $X_j \triangleq \log \frac{\pi_1(Y_{j-1}, Y_j)}{\pi_0(Y_{j-1}, Y_j)}$  and  $\bar{x}_i = E_{\mu_i}[X_1]$  for  $i = 0, 1$ . Derive the analogs of Theorem 2.7.1 and Lemma 2.7.1 by using the results of Section 2.4.3.

## 2.8 Generalized maximum likelihood for finite alphabets

This section is devoted to yet another version of the hypothesis test problem presented in Section 2.7. In particular, the concept of decision test is as in definition 2.7.1 and the associated error probabilities are as given in (2.7.170) there. While the law  $\mu_0$  is again assumed known a-priori, here  $\mu_1$ , the law of  $Y_j$  under the hypothesis  $H_1$  is *unknown*. For that reason, neither the methods nor the results of Section 2.7 apply. Moreover, one has to modify the error criterion since requiring uniformly small  $\beta_n$  over a possibly large class of  $\mu_1$  measures may be too strong (i.e., it may well be that no test can satisfy such a condition). It is reasonable therefore to ask for a criterion which involves asymptotic limits. For finite alphabets  $\Sigma = \{a_1, \dots, a_{|\Sigma|}\}$  such a criterion was suggested by Hoeffding, as follows.

**Definition 2.8.1** *A test  $S$  is optimal (for a given  $\eta > 0$ ) if, among all tests which satisfy*

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log \alpha_n \leq -\eta \quad (2.8.191)$$

*the test  $S$  has maximal exponential rate of error, i.e.  $-\limsup_{n \rightarrow \infty} \{\frac{1}{n} \log \beta_n\}$  is maximal (uniformly over all possible laws  $\mu_1$ ).*

As will become evident in Section ?? a considerable weakening of this criterion is necessary for more general alphabets.

The following lemma states that it suffices to consider functions of the empirical measure when trying to construct an optimal test (i.e., the empirical measure is a sufficient statistic for this problem).

**Lemma 2.8.1** *For any test  $S$  with error probabilities  $\{\alpha_n, \beta_n\}_{n=1}^{\infty}$  there exists a test  $\tilde{S}$  with maps of the form  $\tilde{S}^n(\mathbf{x}) = \tilde{S}(L_n^{\mathbf{x}}, n)$  whose error probabilities  $\{\tilde{\alpha}_n, \tilde{\beta}_n\}_{n=1}^{\infty}$  satisfy*

$$\begin{aligned} \limsup_{n \rightarrow \infty} \frac{1}{n} \log \tilde{\alpha}_n &\leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log \alpha_n \\ \limsup_{n \rightarrow \infty} \frac{1}{n} \log \tilde{\beta}_n &\leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log \beta_n \end{aligned} \quad (2.8.192)$$

**Proof:** For any  $n \in \mathbb{Z}^+$  let  $S_0^n \triangleq (S^n)^{-1}(0)$  and  $S_1^n \triangleq (S^n)^{-1}(1)$  denote the subsets of  $\Sigma^n$  which the maps  $S^n$  assign to  $H_0$  and  $H_1$  respectively. For  $i = 0, 1$  and any  $\nu \in \mathcal{L}_n$  let  $S_i^{\nu, n} \triangleq S_i^n \cap T(\nu)$

(recall that  $T(\nu)$  is the type class of  $\nu$ , see definition 2.1.2). Define

$$\tilde{S}(\nu, n) \triangleq \begin{cases} 0 & \text{if } |S_0^{\nu, n}| \geq \frac{1}{2}|T(\nu)| \\ 1 & \text{otherwise} \end{cases} \quad (2.8.193)$$

where  $|A|$  denotes the cardinality of the set  $A$ . The test  $\tilde{S}$  specified in the lemma's statement is composed of the maps  $\tilde{S}^n(\mathbf{x}) = \tilde{S}(L_n^{\mathbf{x}}, n)$ .

Recall that for any i.i.d. variables  $\mathbf{X} = (X_1, \dots, X_n)$  with marginal  $\mu \in M_1(\Sigma)$  and for any possible type  $\nu \in \mathcal{L}_n$ , the conditional measure  $\text{Prob}_\mu(\mathbf{X} | L_n^{\mathbf{X}} = \nu)$  is a uniform measure over the type class  $T(\nu)$ . In particular, if  $\tilde{S}(\nu, n) = 0$  then

$$\frac{1}{2} \text{Prob}_{\mu_1}(L_n^{\mathbf{X}} = \nu) \leq \frac{|S_0^{\nu, n}|}{|T(\nu)|} \text{Prob}_{\mu_1}(L_n^{\mathbf{X}} = \nu) = \text{Prob}_{\mu_1}(\mathbf{X} \in S_0^{\nu, n}). \quad (2.8.194)$$

Therefore

$$\begin{aligned} \tilde{\beta}_n &= \sum_{\{\nu: \tilde{S}(\nu, n)=0\} \cap \mathcal{L}_n} \text{Prob}_{\mu_1}(L_n^{\mathbf{X}} = \nu) \leq \\ &2 \sum_{\{\nu: \tilde{S}(\nu, n)=0\} \cap \mathcal{L}_n} \text{Prob}_{\mu_1}(\mathbf{X} \in S_0^{\nu, n}) \leq 2 \text{Prob}_{\mu_1}(\mathbf{X} \in S_0^n) = 2\beta_n \end{aligned} \quad (2.8.195)$$

which certainly implies that

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log \tilde{\beta}_n \leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log \beta_n \quad (2.8.196)$$

A similar computation shows that  $\tilde{\alpha}_n \leq 2\alpha_n$ , thus completing the proof.  $\square$

Considering from here on tests which depend only on the empirical type  $L_n^{\mathbf{X}}$ , the following is a characterization of an optimal rule.

**Theorem 2.8.1 (Hoeffding)** *Let the test  $S^*$  consist of the maps*

$$S^{*n}(\mathbf{x}) = \begin{cases} 0 & \text{if } H(L_n^{\mathbf{x}} | \mu_0) < \eta \\ 1 & \text{otherwise} \end{cases} \quad (2.8.197)$$

*Then  $S^*$  is an optimal test.*

**Proof:** By the upper bound of Theorem 2.1.1

$$\begin{aligned} \limsup_{n \rightarrow \infty} \frac{1}{n} \log \text{Prob}_{\mu_0}(H_0 \text{ rejected by } S^*) &= \limsup_{n \rightarrow \infty} \frac{1}{n} \log \text{Prob}_{\mu_0}(L_n^{\mathbf{X}} \in \{\nu : H(\nu | \mu_0) \geq \eta\}) \\ &\leq - \inf_{\{\nu: H(\nu | \mu_0) \geq \eta\}} H(\nu | \mu_0) \leq -\eta \end{aligned} \quad (2.8.198)$$

Therefore,  $S^*$  obviously satisfies the constraint on  $\alpha_n$ . Let  $\beta_n^*$  denote the  $\beta_n$  error probabilities associated with the test  $S^*$ . Then, by the same upper bound (see (2.1.12) in Section 2.1)

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log \beta_n^* = \limsup_{n \rightarrow \infty} \frac{1}{n} \log \text{Prob}_{\mu_1}(L_n^{\mathbf{X}} \in \{\nu : H(\nu|\mu_0) < \eta\}) \leq - \inf_{\{\nu : H(\nu|\mu_0) < \eta\}} H(\nu|\mu_1) \triangleq -I_\eta \quad (2.8.199)$$

where we have used the fact that for (2.1.12) to hold true, the set over which the minimization is performed needs not be closed. When  $I_\eta = \infty$  then  $\limsup_{n \rightarrow \infty} \frac{1}{n} \log \beta_n^* = -\infty$  is the best possible exponential rate. Thus, it suffices to check the optimality of  $S^*$  under laws  $\mu_1$  for which  $I_\eta < \infty$ . Fix one such law  $\mu_1^*$ . Clearly  $\Sigma_{\mu_1^*} \cap \Sigma_{\mu_0}$  is non-empty and moreover

$$\eta^* \triangleq \inf_{\nu \in M_1(\Sigma_{\mu_0} \cap \Sigma_{\mu_1^*})} H(\nu|\mu_0) < \eta, \quad (2.8.200)$$

as  $\max\{H(\nu|\mu_0), H(\nu|\mu_1^*)\} < \infty$  only when  $\Sigma_\nu \subseteq \Sigma_{\mu_0} \cap \Sigma_{\mu_1^*}$ . Furthermore, for any  $\eta' > \eta^*$

$$I_{\eta'} \triangleq \inf_{\{\nu : H(\nu|\mu_0) < \eta'\}} H(\nu|\mu_1^*) = \inf_{\{\nu : H(\nu|\mu_0) < \eta'\} \cap M_1(\Sigma_{\mu_0} \cap \Sigma_{\mu_1^*})} H(\nu|\mu_1^*) < \infty. \quad (2.8.201)$$

Let now  $S$  be any test determined by the binary function  $S(L_n^{\mathbf{X}}, n)$  on  $M_1(\Sigma) \times Z^+$  whose error probabilities  $\alpha_n$  satisfy the constraint (2.8.191). Then, for any  $\delta > 0$  and for all  $n > n_0(\delta)$  large enough

$$\mathcal{L}_n \cap \{\nu : H(\nu|\mu_0) \leq \eta - \delta\} \subset \mathcal{L}_n \cap \{\nu : S(\nu, n) = 0\}. \quad (2.8.202)$$

For otherwise, there exists some  $\delta > 0$  and a sequence of laws  $\nu_n \in \mathcal{L}_n$  (for infinitely many values of  $n$ ), such that  $H(\nu_n|\mu_0) \leq (\eta - \delta)$  while  $S(\nu_n, n) = 1$ . Then, by Lemma 2.1.4, for these values of  $n$ ,

$$\alpha_n \geq \text{Prob}_{\mu_0}(L_n^{\mathbf{X}} = \nu_n) \geq (n+1)^{-|\Sigma|} e^{-nH(\nu_n|\mu_0)} \geq (n+1)^{-|\Sigma|} e^{-n(\eta-\delta)},$$

implying that

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log \alpha_n \geq -(\eta - \delta)$$

and contradicting the constraint (2.8.191). Therefore, by combining (2.8.202) and the lower bound of (2.1.16) one obtains for any  $\delta > 0$

$$\begin{aligned} \liminf_{n \rightarrow \infty} \frac{1}{n} \log \beta_n &\geq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \text{Prob}_{\mu_1^*}(L_n^{\mathbf{X}} \in \{\nu : H(\nu|\mu_0) < \eta - \delta\}) \\ &\geq - \limsup_{n \rightarrow \infty} \left\{ \inf_{\nu \in \mathcal{L}_n : H(\nu|\mu_0) < \eta - \delta} H(\nu|\mu_1^*) \right\} \triangleq -\bar{I}_{\eta-\delta} \end{aligned} \quad (2.8.203)$$

Since for any  $k < \infty$  the set  $\cup_{n=k}^{\infty} \mathcal{L}_n \cap M_1(\Sigma_{\mu_0} \cap \Sigma_{\mu_1^*})$  is a dense subset of  $M_1(\Sigma_{\mu_0} \cap \Sigma_{\mu_1^*})$  over which both  $H(\cdot|\mu_0)$  and  $H(\cdot|\mu_1^*)$  are continuous functions, it follows by (2.8.201) that  $\tilde{I}_{\eta-\delta} = I_{\eta-\delta}$  as long as  $\eta - \delta > \eta^*$ . Thus, one may deduce from (2.8.201) and (2.8.203) that

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log \beta_n \geq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \beta_n \geq - \lim_{\delta \rightarrow 0} \tilde{I}_{\eta-\delta} = - \lim_{\delta \rightarrow 0} I_{\eta-\delta} . \quad (2.8.204)$$

Finally, the optimality of the test  $\mathcal{S}^*$  for the law  $\mu_1^*$  results by comparing (2.8.199) and (2.8.204) provided that

$$\lim_{\delta \rightarrow 0} I_{\eta-\delta} \leq I_{\eta} . \quad (2.8.205)$$

In order to prove (2.8.205) define the measure

$$\mu_0^*(a_i) \triangleq \begin{cases} \mu_0(a_i)/\mu_0(\Sigma_{\mu_1^*}) & \mu_1^*(a_i) > 0 \\ 0 & \text{otherwise} \end{cases} \quad (2.8.206)$$

Note that  $H(\mu_0^*|\mu_0) = -\log \mu_0(\Sigma_{\mu_1^*}) = \eta^*$  (see exercise 2.8.1) and  $\Sigma_{\mu_0^*} = \Sigma_{\mu_0} \cap \Sigma_{\mu_1^*}$ . As  $\{\nu : H(\nu|\mu_0) \leq \eta\}$  is a non-empty compact set there exists a measure  $\nu^*$  such that

$$H(\nu^*|\mu_1^*) \leq I_{\eta} , \quad H(\nu^*|\mu_0) \leq \eta$$

Consider now the family of measures  $\nu_{\theta} = \theta\mu_0^* + (1-\theta)\nu^*$  for  $\theta \in (0, 1)$ . Note that  $\nu_{\theta} \in M_1(\Sigma_{\mu_0} \cap \Sigma_{\mu_1^*})$  and  $\nu_{\theta}$  converges to  $\nu^*$  pointwise as  $\theta \rightarrow 0$  thus implying

$$\lim_{\theta \rightarrow 0} H(\nu_{\theta}|\mu_1^*) = H(\nu^*|\mu_1^*) \leq I_{\eta} . \quad (2.8.207)$$

Moreover, as  $H(\cdot|\mu_0)$  is a convex function

$$H(\nu_{\theta}|\mu_0) \leq \theta H(\mu_0^*|\mu_0) + (1-\theta)H(\nu^*|\mu_0) \leq \theta\eta^* + (1-\theta)\eta$$

implying that for any  $\delta < \theta(\eta - \eta^*)$ ,

$$I_{\eta-\delta} \leq H(\nu_{\theta}|\mu_1^*) \quad (2.8.208)$$

Thus, (2.8.205) follows by combining (2.8.207) and (2.8.208).  $\square$

#### Remarks:

- (a). The finiteness of the alphabet is essential here as (2.8.204) is obtained by applying the lower bounds of Lemma 2.1.4 for individual types instead of the natural large deviations lower

bound for open sets of types. Indeed, for non-finite alphabets a considerable weakening of the optimality criterion is necessary as there are no non-trivial lower bounds for individual types (see Section ??).

- (b). Both Lemma 2.8.1 and Theorem 2.8.1 may be extended to the hypothesis test problem for a known joint law  $\mu_0^n$  versus a family of unknown joint laws  $\mu_1^n$  provided that the random variables  $X_1, \dots, X_n$  are finitely exchangeable under  $\mu_0^n$  and any possible  $\mu_1^n$  so that the empirical measure is still a sufficient statistics. This is outlined in exercises 2.8.5-2.8.6.

### Exercises:

2.8.1 Prove that for any  $\nu \in M_1(\Sigma_{\mu_1^*})$ , if  $\nu \neq \mu_0^*$  then  $H(\nu|\mu_0) > H(\mu_0^*|\mu_0) = \eta^*$ .

2.8.2 Provide an alternative derivation of (2.8.202) based on the results of Section 2.7.

Hint: For any finite alphabet and any known law  $\mu_1$ , deduce from (2.7.171) and exercise 2.7.1 that all probability measures within  $\mathcal{L}_n$  which are of the form  $\nu_\theta(a_i) \triangleq c\mu_1(a_i)^\theta\mu_0(a_i)^{(1-\theta)}$ , where  $\theta \in [0, 1]$  is such that  $H(\nu_\theta|\mu_0) < \eta$  should eventually satisfy  $S(\nu_\theta, n) = 0$ . Then, apply the union of probabilities bound and the volume estimate of Lemma 2.1.1.

2.8.3 (a). Let  $X_j$  for  $j \notin \Delta$  be i.i.d. random variables over the finite alphabet  $\Sigma = \Sigma_{\mu_0}$  while  $X_j$  for  $j \in \Delta$  are unknown deterministic points of  $\Sigma$ . Prove that the test  $S^*$  of (2.8.197) satisfies (2.8.191) for any deterministic increasing sequence of positive integers  $\Delta$  for which  $\lim_{r \rightarrow \infty} \frac{\Delta_r}{r} = \infty$ .

Hint: Let  $L_n^{X^*}$  corresponds to  $\Delta = \emptyset$  and prove that  $\limsup_{n \rightarrow \infty} d_V(L_n^X, L_n^{X^*}) = 0$  almost surely, with some deterministic rate of convergence which depends only upon the sequence  $\Delta$ . Conclude the proof by the continuity of  $H(\cdot|\mu_0)$  over  $M_1(\Sigma)$ .

(b). Construct a counter example to the claim above when  $\Sigma_{\mu_0} \neq \Sigma$ .

2.8.4 Prove that under the assumptions of exercise 2.8.3 part (a), if in addition  $\Sigma_{\mu_1} = \Sigma$  for any possible  $\mu_1$  law then the test  $S^*$  of (2.8.197) is an optimal test. How far can you relax the assumption that  $\Sigma_{\mu_1} = \Sigma$  for all  $\mu_1$  ?



2.8.5 Suppose that for any  $n \in \mathbb{Z}^+$  the random variables  $\mathbf{X}^{(n)} = (X_1, \dots, X_n)$  over the finite alphabet  $\Sigma$  have a known joint law  $\mu_0^n$  under  $H_0$  and an unknown joint law  $\mu_1^n$  under the alternative  $H_1$ . Suppose that for any  $n$  the variables  $\mathbf{X}^{(n)}$  are exchangeable under both  $\mu_0^n$  and any possible  $\mu_1^n$  (namely, the probability of any outcome  $\mathbf{X}^{(n)} = \mathbf{x}$  is invariant under permutations of indices in the vector  $\mathbf{x}$ ). Let  $L_n^{\mathbf{X}}$  denote the empirical measure (type) of  $\mathbf{X}^{(n)}$  and prove that Lemma 2.8.1 holds true in this case. (Note that the variables in  $\mathbf{X}^{(n)}$  may well be dependent).

2.8.6 (a). Consider the scenario described in exercise 2.8.5. Suppose that  $I^n(\nu|\mu_0) = \infty$  implies  $\mu_0^n(L_n^{\mathbf{X}} = \nu) = 0$  and

$$\lim_{n \rightarrow \infty} \sup_{\nu \in \mathcal{L}_n, I^n(\nu|\mu_0) < \infty} \left| \frac{1}{n} \log \mu_0^n(L_n^{\mathbf{X}} = \nu) + I^n(\nu|\mu_0) \right| = 0. \quad (2.8.209)$$

Prove that the test  $\mathcal{S}^*$  of (2.8.197) with  $H(L_n^{\mathbf{X}}|\mu_0)$  replaced by  $I^n(L_n^{\mathbf{X}}|\mu_0)$  is *weakly* optimal in the sense that  $-\limsup_{n \rightarrow \infty} \left\{ \frac{1}{n} \log \beta_n^* \right\}$  is maximal (uniformly over all possible laws  $\mu_1^n$ ) among all the tests for which  $\limsup_{n \rightarrow \infty} \frac{1}{n} \log \alpha_n < \limsup_{n \rightarrow \infty} \frac{1}{n} \log \alpha_n^*$

(b). Apply part (a) to prove the weak optimality of thresholding  $I_{\frac{n}{M}, \mu_M}^{\mathbf{X}}(L_n^{\mathbf{X}})$  of (2.1.25) when testing a given deterministic composition sequence  $\mu_M$  in a sampling without replacement scheme against any unknown composition sequence for such a scheme.

## 2.9 Rate distortion theory for stationary and ergodic sources

Throughout this section we are interested in analyzing the following situation:

be a stationary and ergodic source, with alphabet  $\Sigma$ , i.e.  $\mathcal{P}$  is a stationary ergodic probability measure on  $\Omega = \Sigma^{\mathbb{Z}^+}$ , the space of semi-infinite sequences over  $\Sigma$ . Let  $x_1, x_2, \dots, x_n, \dots$  denote an element of  $\Omega$ , which we say was emitted by the source  $(\mathcal{X}, \mathcal{P})$ . Note that, since  $\mathcal{P}$  is only ergodic, the random variables  $(X_1, \dots, X_n, \dots)$  may well be dependent.

Next, let  $\rho(x, y) : \Sigma \times \Sigma \rightarrow [0, \rho_{\max}]$  be a **one symbol bounded distortion function**, i.e.  $\rho(x, x) = 0$ ,  $\rho(x, y) \neq 0$  for  $x \neq y$  and  $\rho_{\max} < \infty$ . The basic problem of source coding is to find a sequence of deterministic maps (codes)  $C_n : \Sigma^n \rightarrow \Sigma^n$  of small distortion, where the **distortion** of a sequence of codes  $\{C_n\}_{n=1}^{\infty}$  is  $\limsup_{n \rightarrow \infty} \rho_{C_n}$  and

$$\rho_{C_n} \triangleq \frac{1}{n} \sum_{i=1}^n E[\rho(X_i, C_n(\mathbf{X})_i)], \quad (2.9.210)$$

is the average distortion per symbol when the code  $C_n$  is used. The goal of the coding is to allow one to transfer the information contained in sequences emitted by the source  $\mathcal{X}$  with small distortion per symbol, while transmitting as little information as possible.

Clearly, by taking  $C_n$  to be the identity map, one obtains zero distortion but with no coding gain. To get to a more meaningful situation, one would like to have a reduction in the number of possible sequences when using  $C_n$ . Let  $C_n$  also denotes the range of the map  $C_n$  and  $|C_n|$  denotes the cardinality of this set. The **rate** of the code  $C_n$  is defined as

$$R_{C_n} = \frac{1}{n} \log |C_n| \quad (2.9.211)$$

and the smaller  $R_{C_n}$  is, the larger is the coding gain when using  $C_n$ . The main coding theorem, due originally to Shannon, asserts that one cannot hope to get  $R_{C_n}$  to be too small – that indeed, under a bound on the distortion,  $R_{C_n}$  is bounded below in general by some positive quantity and that there are codes which are arbitrarily close to this bound. The proof of this statement which is presented here relies on the large deviations principle of Theorem 2.3.1.

The following definitions are required for the precise statement of the coding theorem.

The distortion associated with any probability measure  $Q$  on  $\Sigma \times \Sigma$  is

$$\rho_Q \triangleq \int_{\Sigma \times \Sigma} \rho(x, y) dQ(x, y). \quad (2.9.212)$$

Let  $Q_X$  and  $Q_Y$  be the marginals of  $Q$ . Then the **mutual information** associated with  $Q$  is

$$H(Q|Q_X \times Q_Y) \triangleq \int_{\Sigma \times \Sigma} \log \left( \frac{dQ(x, y)}{dQ_X(x) dQ_Y(y)} \right) dQ(x, y) \quad (2.9.213)$$

when the above integral is well defined and finite and  $H(Q|Q_X \times Q_Y) = \infty$  otherwise.<sup>1</sup>

The one symbol rate distortion function is defined as

$$R_1(D) = \inf_{\{Q: \rho_Q \leq D, Q_X = P_1\}} H(Q|Q_X \times Q_Y) \quad (2.9.214)$$

where  $P_1$  is the marginal on  $\Sigma$  of the stationary measure  $\mathcal{P}$  on  $\Sigma^{\mathbb{Z}^+}$ .

The one symbol distortion function  $\rho(x, y)$  implies the corresponding  $J$ -symbol average distortion for  $J = 2, 3, \dots$

$$\rho^{(J)}((x_1, \dots, x_J), (y_1, \dots, y_J)) \triangleq \frac{1}{J} \sum_{l=1}^J \rho(x_l, y_l). \quad (2.9.215)$$

---

<sup>1</sup>In information theory books the mutual information is usually denoted by  $I(X; Y)$ . The notation  $H(Q|Q_X \times Q_Y)$  is more consistent with all other notations of this book.

Thus, the  $J$ -symbol distortion associated with a probability measure  $Q$  on  $\Sigma^J \times \Sigma^J$  is

$$\rho_Q^{(J)} = \int_{\Sigma^J \times \Sigma^J} \rho^{(J)}(\mathbf{x}, \mathbf{y}) dQ(\mathbf{x}, \mathbf{y}), \quad (2.9.216)$$

and the *mutual information* associated with the measure  $Q$  having marginals  $Q_X, Q_Y$  on  $\Sigma^J$  is

$$H(Q|Q_X \times Q_Y) \triangleq \int_{\Sigma^J \times \Sigma^J} \log \left( \frac{dQ(\mathbf{x}, \mathbf{y})}{dQ_X(\mathbf{x}) dQ_Y(\mathbf{y})} \right) dQ(\mathbf{x}, \mathbf{y}). \quad (2.9.217)$$

The  $J$ -symbol rate distortion function is therefore defined as

$$R_J(D) \triangleq \inf_{\{Q: \rho_Q^{(J)} \leq D, Q_X = \mathcal{P}_J\}} \frac{1}{J} H(Q|Q_X \times Q_Y), \quad (2.9.218)$$

where  $\mathcal{P}_J$  is the  $J$ -th marginal (on  $\Sigma^J$ ) of the stationary measure  $\mathcal{P}$ . Finally, the **rate distortion function** is

$$R(D) \triangleq \inf_{J \geq 1} R_J(D). \quad (2.9.219)$$

The source coding theorem states that the rate distortion function is a tight lower bound on the limiting rate  $R_{C_n}$  of a sequence of codes  $\{C_n\}_{n=1}^\infty$  with distortion  $D$ .

#### Theorem 2.9.1 (Source Coding Theorem)

- (a). **Direct Part:** For any  $D \geq 0$  such that  $R(D) < \infty$  and any  $\delta > 0$ , there exists a sequence of codes  $\{C_n\}_{n=1}^\infty$  with distortion at most  $D$  and rates  $R_{C_n} \leq R(D) + \delta$ .
- (b). **Converse Part:** For any sequence of codes  $\{C_n\}_{n=1}^\infty$  of distortion  $D$  and any  $\delta > 0$   $\liminf_{n \rightarrow \infty} R_{C_n} \geq R(D + \delta)$ .

**Remark:** Note that  $|\Sigma|$  may be infinite and there are no structural conditions on  $\Sigma$  besides the requirement that  $\mathcal{P}$  be based on  $\Sigma^{\mathbb{Z}^+}$ . On the other hand, whenever  $R(D)$  is finite, the resulting codes always take values in some finite set and in particular, may be represented by finite binary sequences.

The proof of the Source Coding Theorem is presented via a sequence of lemmas with the large deviations principle of Section 2.3 implying the first lemma which is key for the Direct Part of the theorem.

**Lemma 2.9.1** Suppose  $Q$  is any probability measure on  $\Sigma \times \Sigma$  for which  $I_Q \triangleq H(Q|Q_X \times Q_Y) < \infty$ ,  $Q_X \equiv \mathcal{P}_1$  and  $\rho_Q < \rho_{Q_X \times Q_Y}$  (where  $\rho_{Q_X \times Q_Y}$  is the distortion associated with  $Q_X \times Q_Y \equiv \mathcal{P}_1 \times$

$Q_Y$ ). Let  $Z_n(\mathbf{x}) \triangleq \frac{1}{n} \sum_{j=1}^n \rho(x_j, Y_j)$  where  $Y_j$  are independent random variables distributed over  $\Sigma$  according to  $Q_Y$ . Then,

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log \text{Prob}(Z_n(\mathbf{X}) < \rho_Q | \mathbf{X}) \geq -I_Q \quad \text{almost surely } \mathcal{P} \quad (2.9.220)$$

where  $\mathbf{X}$  is a random sequence of symbols emitted by the source  $\mathcal{X}$ .

**Proof:** The inequality (2.9.220) follows from the lower bound of Theorem 2.3.1 for the open set  $\Gamma = (-\infty, \rho_Q)$  and the random variables  $\{Z_n\}_{n=1}^\infty$ . Specifically, Theorem 2.3.1 is applied per element of  $\Omega = \Sigma^{Z^+}$ , where the conditions of this theorem hold almost surely by the ergodicity of  $\mathcal{P}$ . To verify that, let  $\Lambda_n(\theta) \triangleq \log E[e^{\theta Z_n(\mathbf{X})} | \mathbf{X}]$  and note that since  $\rho(\cdot, \cdot)$  is a bounded function  $|\Lambda_n(\theta)| < \infty$  for all  $\theta \in \mathbb{R}$  and any realization of  $\mathbf{X}$ . By Birkhoff's ergodic theorem [3]

$$\Lambda(\theta) \triangleq \lim_{n \rightarrow \infty} \frac{1}{n} \Lambda_n(n\theta) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n \log \int_{\Sigma} e^{\theta \rho(X_j, y)} dQ_Y(y) \quad (2.9.221)$$

exists almost surely  $\mathcal{P}$ . Moreover,

$$\Lambda(\theta) = \int_{\Sigma} \log \left( \int_{\Sigma} e^{\theta \rho(x, y)} dQ_Y(y) \right) d\mathcal{P}_1(x) = \int_{\Sigma} \log \left( \int_{\Sigma} e^{\theta \rho(x, y)} dQ_Y(y) \right) dQ_X(x), \quad (2.9.222)$$

does not depend on the specific sequence  $\mathbf{X}$  emitted by the source (recall that  $Q_X \equiv \mathcal{P}_1$  implying the second equality above). Furthermore, since  $\rho(\cdot, \cdot)$  is uniformly bounded the function  $\Lambda(\cdot)$  is finite and differentiable everywhere (in  $\mathbb{R}^1$ ). Therefore, by Lemma 2.3.2, part (c) of Theorem 2.3.1 applies, yielding

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log P(Z_n(\mathbf{X}) < \rho_Q | \mathbf{X}) \geq - \inf_{x < \rho_Q} \Lambda^*(x) \triangleq -J_Q \quad \text{almost surely } \mathcal{P} \quad (2.9.223)$$

Recall that  $\Lambda^*(x) \triangleq \sup_{\lambda} (\lambda x - \Lambda(\lambda))$ . As  $\Lambda(0) = 0$  (see (2.9.222)),  $\Lambda'(0) = \rho_{Q_X \times Q_Y}$  (differentiate (2.9.222) and compare with (2.9.212)) and  $\Lambda'(\lambda)$  is monotonically nondecreasing (since  $\rho(\cdot, \cdot) \geq 0$ ) it follows that for  $x \leq \rho_Q < \Lambda'(0)$

$$\sup_{\lambda > 0} [\lambda x - \Lambda(\lambda)] \leq \sup_{\lambda > 0} [\lambda \Lambda'(0) - \Lambda(\lambda)] \leq \Lambda^*(\Lambda'(0)) = 0 \Lambda'(0) - \Lambda(0) = 0. \quad (2.9.224)$$

Since  $\Lambda^*(x) \geq 0$ , it follows that for  $x \leq \rho_Q$

$$\Lambda^*(x) = \sup_{\lambda \leq 0} [\lambda x - \Lambda(\lambda)] \geq \Lambda^*(\rho_Q) \quad (2.9.225)$$

and thus  $J_Q = \Lambda^*(\rho_Q)$ .

For any  $\lambda \in \mathbb{R}^1$  define the probability measure  $Q_\lambda$  on  $\Sigma \times \Sigma$  via

$$\frac{dQ_\lambda(x, y)}{dQ_X(x)dQ_Y(y)} \triangleq \frac{e^{\lambda \rho(x, y)}}{\int_\Sigma e^{\lambda \rho(x, z)} dQ_Y(z)}.$$

Since  $\rho(\cdot, \cdot)$  is a bounded function the measures  $Q_\lambda$  and  $Q_X \times Q_Y$  are equivalent and as  $H(Q|Q_X \times Q_Y) < \infty$  the relative entropy  $H(Q|Q_\lambda)$  is well defined. Moreover,

$$\begin{aligned} 0 &\leq H(Q|Q_\lambda) \triangleq \int_{\Sigma \times \Sigma} \log \frac{dQ(x, y)}{dQ_\lambda(x, y)} dQ(x, y) = \\ &= \int_{\Sigma \times \Sigma} \log \left\{ \frac{dQ(x, y) e^{-\lambda \rho(x, y)}}{dQ_X(x) dQ_Y(y)} \int_\Sigma e^{\lambda \rho(x, z)} dQ_Y(z) \right\} dQ(x, y) = I_Q - \lambda \rho_Q + \Lambda(\lambda) \end{aligned} \quad (2.9.226)$$

Since (2.9.226) holds for all  $\lambda$ , it follows that  $I_Q \geq \Lambda^*(\rho_Q)$ . The proof is now completed in view of (2.9.223) and (2.9.225).  $\square$

The proof of the Direct Part of Theorem 2.9.1 is based on a random coding argument, where instead of specifically constructing the codes  $C_n$ , the classes  $C_n$  of all codes of some fixed size are considered. Let  $\bar{\rho}_n \triangleq E_{C_n}[\rho_{C_n}]$  be the average of  $\rho_{C_n}$  over  $C_n$ , where the distribution within the class  $C_n$  results by choosing the code  $C_n$  at random according to some probability measure. For any probability measure over  $C_n$  there exists at least one code in  $C_n$  for which  $\rho_{C_n} \leq \bar{\rho}_n$ . In the following lemma an upper bound on  $\bar{\rho}_n$  is derived based on the large deviations lower bound of Lemma 2.9.1.

**Lemma 2.9.2** Suppose  $Q$  is a probability measure on  $\Sigma \times \Sigma$  for which  $H(Q|Q_X \times Q_Y) < \infty$  and  $Q_X \equiv \mathcal{P}_1$ . Fix  $\delta > 0$  arbitrarily small and let  $C_n$  be the class of all codes  $C_n$  of size  $|C_n| = \lfloor e^{n(H(Q|Q_X \times Q_Y) + \delta)} \rfloor$ . Then, there exist distributions on  $C_n$  for which  $\limsup_{n \rightarrow \infty} \bar{\rho}_n \leq \rho_Q$ .

**Proof:** Fix the probability measure  $Q$  and let  $I_Q \triangleq H(Q|Q_X \times Q_Y)$ .

(a). Suppose that  $\rho_{Q_X \times Q_Y} \leq \rho_Q$  (recall that  $\rho_{Q_X \times Q_Y}$  is the distortion associated with the measure  $Q_X \times Q_Y \equiv \mathcal{P}_1 \times Q_Y$ ). Let  $Y_1, \dots, Y_n$  be i.i.d according to the law  $Q_Y$  and  $\mathbf{Y} \triangleq (Y_1, \dots, Y_n)$  be a code-word of  $C_n$  to which all  $\Sigma^n$  is mapped. Since  $I_Q \geq 0$  this construction is always possible ( $|C_n| \geq 1$ ) and it results with a class of codes  $C_n$ , one code per realization of  $\mathbf{Y}$ . The average distortion over codes in  $C_n$  is now

$$\bar{\rho}_n = E_{\mathbf{Y}} \left[ \frac{1}{n} \sum_{j=1}^n E[\rho(X_j, Y_j) | \mathbf{Y}] \right] = \int_{\Sigma \times \Sigma} \rho(x, y) d\mathcal{P}_1(x) dQ_Y(y) = \rho_{Q_X \times Q_Y}. \quad (2.9.227)$$

Thus, since  $\bar{\rho}_n = \rho_{Q_X \times Q_Y} \leq \rho_Q$  the proof of the lemma is complete.

(b). Consider now the case when  $\rho_{Q_X \times Q_Y} > \rho_Q$ ,  $I_Q < \infty$  and  $Q_X \equiv \mathcal{P}_1$ . Let

$k_n \triangleq |C_n| = \lfloor e^{n(I_Q + \delta)} \rfloor$  be as specified in the statement of the lemma and  $Y_j^{(i)}$ ,  $j = 1, \dots, n$ ,  $i = 1, \dots, k_n$ , be  $n \times k_n$  i.i.d. random variables of law  $Q_Y$ . The probability distribution on the class of codes  $C_n$  is generated by considering the codes with code-words  $\mathbf{Y}^{(i)} \triangleq (Y_1^{(i)}, \dots, Y_n^{(i)})$  for  $i = 1, \dots, k_n$ . Per realization of these code-words the mapping  $C_n$  is constructed as follows. For any  $\mathbf{x} \in \Sigma^n$  define the set

$$S_n(\mathbf{x}) \triangleq \{y : \frac{1}{n} \sum_{j=1}^n \rho(x_j, y_j) < \rho_Q\} \quad (2.9.228)$$

and let  $C_n(\mathbf{x})$  be any element of  $C_n \cap S_n(\mathbf{x})$ , where if this set is empty then  $C_n(\mathbf{x})$  is arbitrarily chosen. For this mapping,

$$\frac{1}{n} \sum_{j=1}^n \rho(x_j, C_n(\mathbf{x})_j) \leq \rho_Q + \rho_{\max} 1_{C_n \cap S_n(\mathbf{x}) = \emptyset}$$

implying that

$$\bar{\rho}_n \leq \rho_Q + \rho_{\max} \text{Prob}(C_n \cap S_n(\mathbf{X}) = \emptyset), \quad (2.9.229)$$

where  $\mathbf{X} \in \Sigma^n$  is a random sequence of  $n$  symbols emitted by the source  $\mathcal{X}$  and the set  $C_n$  consists of the  $k_n$  i.i.d. random vectors  $\mathbf{Y}^{(i)}$ . Clearly,

$$\begin{aligned} \text{Prob}(C_n \cap S_n(\mathbf{X}) = \emptyset) &= E[\text{Prob}(\mathbf{Y}^{(i)} \notin S_n(\mathbf{X}) \text{ for all } i | \mathbf{X})] \\ &= E[(1 - \text{Prob}(\mathbf{Y}^{(1)} \in S_n(\mathbf{X}) | \mathbf{X}))^{k_n}] \leq E[e^{-k_n \text{Prob}(\mathbf{Y}^{(1)} \in S_n(\mathbf{X}) | \mathbf{X})}]. \end{aligned} \quad (2.9.230)$$

By (2.9.229) and (2.9.230),  $\limsup_{n \rightarrow \infty} \bar{\rho}_n \leq \rho_Q$  provided that  $k_n \text{Prob}(\mathbf{Y}^{(1)} \in S_n(\mathbf{X}) | \mathbf{X}) \rightarrow \infty$  as  $n \rightarrow \infty$ , in probability  $\mathcal{P}$ . By the definition of  $k_n$  it suffices to show that

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log \text{Prob}(\mathbf{Y}^{(1)} \in S_n(\mathbf{X}) | \mathbf{X}) \geq -I_Q \quad \text{in Probability } \mathcal{P}, \quad (2.9.231)$$

in order to complete the proof of the lemma. Since Lemma 2.9.1 applies to  $Z_n(\mathbf{x}) \triangleq \frac{1}{n} \sum_{j=1}^n \rho(x_j, Y_j^{(1)})$  and  $\{\mathbf{Y}^{(1)} \in S_n(\mathbf{x})\} \equiv \{Z_n(\mathbf{x}) < \rho_Q\}$  the bound (2.9.231) follows.  $\square$

The following weak version of the direct part of the Source Coding Theorem is an immediate consequence of Lemma 2.9.2.

**Lemma 2.9.3** *For any  $D \geq 0$  such that  $R_1(D) < \infty$  and any  $\delta > 0$  there exists a sequence of codes  $\{C_n\}_{n=1}^{\infty}$  with distortion at most  $D$  and rates  $R_{C_n} \leq R_1(D) + \delta$ .*

**Proof:** Since  $R_1(D) < \infty$  there exists a sequence of measures  $\{Q^{(m)}\}_{m=1}^{\infty}$  such that  $I_{Q^{(m)}} \rightarrow R_1(D)$  while  $\rho_{Q^{(m)}} \leq D$  and  $Q_X^{(m)} \equiv \mathcal{P}_1$ . By applying Lemma 2.9.2 for  $Q^{(m)}$ ,  $m = 1, 2, \dots$ , it follows that  $\limsup_{n \rightarrow \infty} \bar{\rho}_n \leq \limsup_{m \rightarrow \infty} \rho_{Q^{(m)}} \leq D$  when  $C_n$  is the class of all codes of size  $\lfloor e^{n(R_1(D)+\delta)} \rfloor$  and  $\delta > 0$  is arbitrarily small. The existence of a sequence of codes  $C_n$  of rates  $R_{C_n} \leq R_1(D) + \delta$  and of distortion  $\limsup_{n \rightarrow \infty} \rho_{C_n} \leq D$  is deduced by extracting the codes  $C_n$  of minimal  $\rho_{C_n}$  from the ensembles  $C_n$ .  $\square$

For i.i.d. source symbols,  $R(D) = R_1(D)$  (see exercise 2.9.2) and the Direct Part of the Source Coding Theorem amounts to Lemma 2.9.3. When the symbols are dependent one needs the following extension of Lemma 2.9.2.

**Lemma 2.9.4** *Suppose  $\mathcal{P}$  is ergodic with respect to the  $J$ -th shift operation (namely, it is ergodic in blocks of size  $J$ ). Then, for any probability measure  $Q$  on  $\Sigma^J \times \Sigma^J$  with  $Q_X \equiv \mathcal{P}_J$  and for any  $\delta > 0$  there exists a sequence of codes  $C_n$  of rates  $R_{C_n} \leq \frac{1}{J}H(Q|Q_X \times Q_Y) + \delta$  and of distortion at most  $\rho_Q^{(J)}$ .*

**Proof:** Consider the enlarged alphabet  $\Sigma^J$  and regard each block of  $J$  consecutive symbols of the emitted source sequence as one symbol from  $\Sigma^J$ . A sketch of the proof which basically follows the proofs of Lemmas 2.9.1 and 2.9.2 is presented here. Let

$$\Lambda^{(J)}(\theta) = \frac{1}{J} \int_{\Sigma^J} \log \left( \int_{\Sigma^J} e^{J\theta \rho^{(J)}(x,y)} dQ_Y(y) \right) d\mathcal{P}_J(x). \quad (2.9.232)$$

Then, by the ergodicity of  $\mathcal{P}$  in blocks of size  $J$

$$\Lambda^{(J)}(\theta) = \lim_{n \rightarrow \infty} \frac{1}{nJ} \log E[e^{n\theta J Z_{nJ}^{(J)}(\mathbf{X})} | \mathbf{X}] \text{ almost surely } \mathcal{P}, \quad (2.9.233)$$

where  $Z_{nJ}^{(J)}(\mathbf{x}) \triangleq \frac{1}{n} \sum_{j=1}^n \rho^{(J)}(\mathbf{x}_j, \mathbf{Y}_j)$  and  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  are concatenated symbols namely, are elements of  $\Sigma^J$  (while  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$  are i.i.d. random variables of law  $Q_Y$  over  $\Sigma^J$ ). Thus, Theorem 2.3.1 is once again applicable. To complete the proof, define  $Q_\lambda$  now via

$$\frac{dQ_\lambda(\mathbf{x}, \mathbf{y})}{dQ_X(\mathbf{x})dQ_Y(\mathbf{y})} \triangleq \frac{e^{J\lambda \rho^{(J)}(\mathbf{x}, \mathbf{y})}}{\int_{\Sigma^J} e^{J\lambda \rho^{(J)}(\mathbf{x}, \mathbf{z})} dQ_Y(\mathbf{z})},$$

and observe that for any  $\lambda \in \mathbb{R}^1$  and any probability measure  $Q$  over  $\Sigma^J \times \Sigma^J$  with  $Q_X \equiv \mathcal{P}_J$

$$0 \leq \frac{1}{J}H(Q|Q_\lambda) = \frac{1}{J}H(Q|Q_X \times Q_Y) - \lambda \rho_Q^{(J)} + \Lambda^{(J)}(\lambda).$$

Thus,  $\frac{1}{J}H(Q|Q_X \times Q_Y) \geq \Lambda^{(J)*}(\rho_Q^{(J)})$ . □

The following corollary is an immediate consequence of Lemma 2.9.4 (by adapting the proof of Lemma 2.9.3 to  $J \neq 1$ ).

**Corollary 2.9.1** *If  $\mathcal{P}$  is ergodic in blocks of size  $J$  for any  $J \in \mathbb{Z}^+$  then the direct part of the Source Coding Theorem holds.*

While in general, an ergodic  $\mathcal{P}$  might be non ergodic in blocks, Corollary 2.9.1 holds true for any stationary and ergodic  $\mathcal{P}$  (see for example [2], pp. 278–280, [13], pp. 496–500). It clearly suffices for that purpose to prove the following lemma.

**Lemma 2.9.5** *Lemma 2.9.4 holds true for any stationary and ergodic  $\mathcal{P}$  and any  $J \in \mathbb{Z}^+$ .*

**Proof:** It is possible to show that when considering blocks of size  $J$ , the emitted infinite sequences of the source may almost surely be divided into  $J$  equally probable ergodic modes,  $E_0, \dots, E_{J-1}$ , such that if sequence  $(x_1, x_2, \dots, x_n, \dots)$  belongs to mode  $E_i$  then  $(x_{1+k}, x_{2+k}, \dots, x_{n+k}, \dots)$  belongs to the mode  $E_{(i+k) \bmod J}$  (see [2],[13]). This implies that  $\mathcal{P} = \frac{1}{J} \sum_{i=0}^{J-1} \mathcal{P}^{(i)}$  where  $\{\mathcal{P}^{(i)}\}_{i=0}^{J-1}$  are ergodic in blocks of size  $J$  and correspond to the ergodic modes  $E_i$ . By projections onto these ergodic modes, each law  $Q$  on  $\Sigma^J \times \Sigma^J$  with  $Q_X \equiv \mathcal{P}_J$  is similarly decomposed into  $Q = \frac{1}{J} \sum_{i=0}^{J-1} Q^{(i)}$  such that  $Q_X^{(i)} = \mathcal{P}_J^{(i)}$  is the  $J$ -th marginal of  $\mathcal{P}^{(i)}$  while  $\rho_Q^{(J)} = \frac{1}{J} \sum_{i=0}^{J-1} \rho_{Q^{(i)}}^{(J)}$  and

$$H(Q|Q_X \times Q_Y) \geq \frac{1}{J} \sum_{i=0}^{J-1} H(Q^{(i)}|Q_X^{(i)} \times Q_Y^{(i)}). \quad (2.9.234)$$

By applying Lemma 2.9.4 for  $\{Q^{(i)}\}_{i=0}^{J-1}$  there exist  $J$  sequences of codes  $C_n^{(i)}$ ,  $i = 0, \dots, J-1$  with rates  $R_{C_n^{(i)}} \leq \frac{1}{J}H(Q^{(i)}|Q_X^{(i)} \times Q_Y^{(i)}) + \delta$  and distortions at most  $\rho_{Q^{(i)}}^{(J)}$  with respect to the source measures  $\{\mathcal{P}^{(i)}\}_{i=0}^{J-1}$ , respectively. A sequence of codes  $C_{(n+1)J}^*$  of rates

$$R_{C_{(n+1)J}^*} \leq \frac{1}{J} \sum_{i=0}^{J-1} R_{C_n^{(i)}} + \frac{1}{(n+1)J} \log J \quad (2.9.235)$$

and of distortion at most  $\rho_Q^{(J)}$  with respect to  $\mathcal{P}$  is now constructed as follows. The code  $C_{(n+1)J}^*$  is the union of  $J$  codes  $\{\tilde{C}_{(n+1)J}^{(i)}\}_{i=0}^{J-1}$ , each of cardinality  $\prod_{i=0}^{J-1} |C_n^{(i)}|$  and length  $(n+1)J$ . The code  $\tilde{C}_{(n+1)J}^{(i)}$  consists of all the distinct words in



$\{(y_0, a^*, y_1, a^*, \dots, y_{J-1}, a^*) : y_k \in C_n^{(i+k) \bmod J}, k = 0, \dots, J-1\}$  and  $a^* \in \Sigma$  is a fixed separator symbol. This construction guarantees that the sequence of codes  $\tilde{C}_{(n+1)J}^{(i)}$  has distortion at most  $\rho_Q^{(J)}$  with respect to the source measure  $\mathcal{P}^{(i)}$  and thus  $C_{(n+1)J}^*$  has at most this distortion with respect to the source measure  $\mathcal{P}$ . By (2.9.234) and (2.9.235) for all  $n$  large enough,  $R_{C_{(n+1)J}^*} \leq \frac{1}{J}H(Q|Q_X \times Q_Y) + 2\delta$ . Finally, for code length which is not an integer multiple of  $J$ , one may easily modify the code  $C_{(n+1)J}^*$  of closest length while neither affecting the *limiting* distortion per symbol nor the limiting rate (as  $n \rightarrow \infty$ ). The proof is thus complete.  $\square$

The last lemma in this section is devoted to the converse part of the Source Coding Theorem, whose proof is based on information theoretical arguments and not on large deviations bounds.

**Lemma 2.9.6** *For any sequence of codes  $\{C_n\}_{n=1}^\infty$  of distortion  $D$  and for any  $\delta > 0$*   
 $\liminf_{n \rightarrow \infty} R_{C_n} \geq R(D + \delta)$ .

**Proof:** It suffices to consider codes  $C_n$  of finite rates and of distortion  $D$ . Such a code  $C_n$  is a mapping from  $\Sigma^n$  to  $\Sigma^n$ . When its domain  $\Sigma^n$  is equipped with the probability measure  $\mathcal{P}_n$  (the  $n$ -th marginal of the source measure  $\mathcal{P}$ ),  $C_n$  induces a (degenerate) joint measure  $Q^{(n)}$  on  $\Sigma^n \times \Sigma^n$ . Note that  $Q_X^{(n)} \equiv \mathcal{P}_n$  and  $\rho_{Q^{(n)}} = \rho_{C_n} \leq D + \delta$  for all  $n$  large enough (and any  $\delta > 0$ ). Therefore,  $H(Q^{(n)}|Q_X^{(n)} \times Q_Y^{(n)}) \geq nR_n(D + \delta) \geq nR(D + \delta)$  for any  $\delta > 0$  and any  $n$  large enough. Since the marginals  $Q_Y^{(n)}$  have the finite support sets  $C_n$

$$nR_{C_n} = \log |C_n| \geq H(Q_Y^{(n)}),$$

where the entropy  $H(Q_Y^{(n)})$  is defined as  $H(Q_Y^{(n)}) \triangleq \sum_{i=1}^{|C_n|} Q_Y^{(n)}(y_i) \log \frac{1}{Q_Y^{(n)}(y_i)}$  (compare with the definition in Section 2.1.1). Let  $f_n(x, y_i) \triangleq \frac{dQ^{(n)}(x, y_i)}{dQ_X^{(n)}(x) dQ_Y^{(n)}(y_i)}$ . Then,  $f_n : \Sigma^n \times \Sigma^n \rightarrow [0, \infty)$  is well defined,  $\int_{\Sigma^n} f_n(x, y_i) dQ_X^{(n)}(x) = 1$  for all  $i$  as well as  $\sum_{i=1}^{|C_n|} f_n(x, y_i) Q_Y^{(n)}(y_i) = 1$  almost surely  $Q_X^{(n)}$ . Thus,  $f_n(x, y_i) Q_Y^{(n)}(y_i) \leq 1$  almost surely  $Q_X^{(n)}$  and

$$\begin{aligned} H(Q_Y^{(n)}) - H(Q^{(n)}|Q_X^{(n)} \times Q_Y^{(n)}) &= \sum_{i=1}^{|C_n|} Q_Y^{(n)}(y_i) \left[ \log \frac{1}{Q_Y^{(n)}(y_i)} - \int_{\Sigma^n} f_n(x, y_i) \log f_n(x, y_i) dQ_X^{(n)}(x) \right] \\ &= \int_{\Sigma^n} dQ_X^{(n)}(x) \left\{ \sum_{i=1}^{|C_n|} f_n(x, y_i) Q_Y^{(n)}(y_i) \log \left[ \frac{1}{f_n(x, y_i) Q_Y^{(n)}(y_i)} \right] \right\} \geq 0 \end{aligned} \quad (2.9.236)$$

Therefore,

$$\liminf_{n \rightarrow \infty} R_{C_n} \geq \liminf_{n \rightarrow \infty} \frac{1}{n} H(Q_Y^{(n)}) \geq \liminf_{n \rightarrow \infty} \frac{1}{n} H(Q^{(n)}|Q_X^{(n)} \times Q_Y^{(n)}) \geq R(D + \delta).$$

□

**Exercises:**

**2.9.1** (a). Prove that when  $\Sigma$  is a finite set and  $R_1(D) > 0$ , there exists a probability measure  $Q$  on  $\Sigma \times \Sigma$  for which  $I_Q = R_1(D)$ ,  $Q_X = \mathcal{P}_1$  and  $\rho_Q = D$ .

(b). Prove that for this measure also  $\Lambda^*(\rho_Q) = I_Q$ .

**2.9.2** Prove that when  $\mathcal{P}$  is a product measure (namely, the emitted symbols  $X_1, X_2, \dots, X_n$  are i.i.d.) then  $R_J(D) = R_1(D)$  for all  $J \in \mathbb{Z}^+$ .

**2.9.3** (a). Show that  $(m+n)R_{m+n}(D) \leq mR_m(D) + nR_n(D)$  for any two integers  $m, n$ .

(b). Conclude that if  $\limsup_{J \rightarrow \infty} R_J(D) < \infty$  then  $R(D) = \lim_{J \rightarrow \infty} R_J(D)$ .

emphasize in intr. that deal with iid case. correctiuons are denoted

## 2.10 Refinements of large deviations statements in $\mathbb{R}^d$

Cramer's theorem deals with the tails of the empirical means  $\hat{S}_n$ . On a finer scale, at least in the i.i.d. case, the random variables  $\sqrt{n}\hat{S}_n$  possess a limiting Normal distribution by the classical Central Limit Theorem. In this situation, the empirical means  $n^\beta \hat{S}_n$  satisfy a large deviations principle for any  $\beta \in (0, \frac{1}{2})$ , but always with a quadratic (Normal-like) rate function. This statement is made precise in the following theorem.

**Theorem 2.10.1** *Let  $X_1, \dots, X_n$  be a sequence of  $\mathbb{R}^d$  valued i.i.d. random variables with  $E(X_i) = 0$  such that  $\Lambda_X(\lambda) \triangleq \log E[e^{\langle \lambda, X_i \rangle}] < \infty$  in some open ball  $B_{0,\delta}$  around the origin. Fix  $\beta \in (0, \frac{1}{2})$  and let  $Z_n \triangleq \frac{1}{n^{1-\beta}} \sum_{i=1}^n X_i = n^\beta \hat{S}_n$ . Then,  $Z_n$  satisfy a large deviations principle in  $\mathbb{R}^d$  governed by the good rate function*

$$I_g(x) \triangleq \frac{1}{2} \langle x, Cx \rangle, \quad (2.10.237)$$

where  $C$  is the covariance matrix of  $X_i$ . Moreover, any open or closed set  $G$  is an  $I_g$  continuity set.

**Proof:** This theorem follows from the general large deviations statement of Section 2.3 with  $a_n \triangleq n^{(2\beta-1)}$ . Indeed, in the notations of Section 2.3

$$\begin{aligned}\Lambda_n(a_n^{-1}\lambda) &= \log E\left(e^{a_n^{-1}\langle\lambda, Z_n\rangle}\right) \\ &= \sum_{i=1}^n \log E\left(e^{n^{-\beta}\langle\lambda, X_i\rangle}\right) = n \log E\left(e^{n^{-\beta}\langle\lambda, X_1\rangle}\right)\end{aligned}\quad (2.10.238)$$

Therefore

$$\Lambda(\lambda) = \lim_{n \rightarrow \infty} n^{2\beta} \log E\left(e^{n^{-\beta}\langle\lambda, X_1\rangle}\right) \quad (2.10.239)$$

Note that  $n^{-\beta} \xrightarrow[n \rightarrow \infty]{} 0$  and therefore, by our assumption that  $\Lambda(\lambda) < \infty$  in an open ball around 0, for each  $\lambda \in \mathbb{R}^d$  there exists an  $n_0$  large enough such that for all  $n > n_0$ ,  $E\left(e^{n^{-\beta}\langle\lambda, X_1\rangle}\right) < \infty$ , and by dominated convergence

$$E\left(e^{n^{-\beta}\langle\lambda, X_1\rangle}\right) = 1 + n^{-\beta} E[\langle\lambda, X_1\rangle] + \frac{1}{2} n^{-2\beta} E[\langle\lambda, X_1\rangle^2] + O(n^{-3\beta}) \quad (2.10.240)$$

Substituting (2.10.240) into (2.10.239) one obtains (using the identity  $E[\langle\lambda, X_1\rangle] = 0$ )

$$\begin{aligned}\Lambda(\lambda) &= \lim_{n \rightarrow \infty} n^{2\beta} \log \left\{ 1 + \frac{1}{2} n^{-2\beta} E[\langle\lambda, X_1\rangle^2] + O(n^{-3\beta}) \right\} \\ &= \frac{1}{2} E[\langle\lambda, X_1\rangle^2] = \frac{1}{2} \langle\lambda, C\lambda\rangle\end{aligned}\quad (2.10.241)$$

Thus,

$$\Lambda^*(x) \triangleq \sup_{\lambda \in \mathbb{R}^d} \{ \langle\lambda, x\rangle - \Lambda(\lambda) \} = \sup_{\lambda \in \mathbb{R}^d} \{ \langle\lambda, x\rangle - \frac{1}{2} \langle\lambda, C\lambda\rangle \} = \frac{1}{2} \langle x, Cx \rangle = I_g(x) \quad (2.10.242)$$

Since  $\Lambda^*(\lambda) = \frac{1}{2} \langle\lambda, C\lambda\rangle$  is differentiable and finite everywhere Theorem 2.3.1 applies and the proof is thus complete.  $\square$

#### Remarks:

(a). Note that Theorem 2.10.1 is nothing more than what one would obtain from a naive Taylor expansion applied on the anastaz  $\text{Prob}(\hat{S}_n = x) \approx e^{-nI(x)}$  where  $I(\cdot)$  is the rate function of Theorem 2.3.2 (see Section 2.3).

(b). The rate of convergence in (2.10.241) is of  $O(n^{-\beta})$ , suggesting a similar convergence rate for  $a_n \log \text{Prob}(Z_n \in G)$  (which converges to  $-\inf_{x \in G} I_g(x)$  for any open  $G$ ). Indeed, such a result is proved in [12] (pp. 552-553) for  $d = 1$  and  $G = (x, \infty)$ . This extends for example the validity of

the Normal approximation for the distribution of  $\sqrt{n}\hat{S}_n$  to intervals of  $o(n^{\frac{1}{6}})$  (which correspond to  $\beta > \frac{1}{3}$  here).

(c). A similar result may be obtained in the context of Markov additive processes (see Section 2.4.1).

Another refinement of Cramer's theorem involves a more accurate estimate of the laws  $\mu_n$  of  $\hat{S}_n$  (the empirical means of i.i.d. random variables). Specifically, for a "nice"  $I$  Continuity Set  $A$  one seeks an estimate  $J_n$  of  $\mu_n(A)$  such that  $\lim_{n \rightarrow \infty} J_n \mu_n(A) = 1$ . Such an estimate is an improvement over the normalized logarithmic limit  $\frac{1}{n} \log \mu_n(A)$  implied by a large deviations principle. The following theorem deals with the estimate  $J_n$  for certain half intervals  $A = [q, \infty) \subset \mathbb{R}^1$ .

**Theorem 2.10.2 (Bahadur and Rao)** *Let  $\mu_n$  denotes the law of  $\hat{S}_n = \frac{1}{n} \sum_{i=1}^n X_i$  where  $X_i$  are i.i.d. real valued random variables and  $\Lambda(\lambda) = \log E[e^{\lambda X_1}]$  is the logarithmic moment generating function of  $X_i$ . Consider the set  $A = [q, \infty)$  where  $q \in \mathcal{F}$ , namely  $q = \Lambda'(\eta)$  for some positive  $\eta \in \mathcal{D}_\Lambda^0$ .*

(a). *If the law of  $X_1$  is non-lattice, then*

$$\lim_{n \rightarrow \infty} J_n \mu_n(A) = 1 \quad (2.10.243)$$

where  $J_n = \eta \sqrt{\Lambda''(\eta)} 2\pi n e^{n\Lambda^*(q)}$ .

(b). *Suppose  $X_1$  has a lattice law, namely, for some finite  $x_0, d$ , the random variable  $\frac{1}{d}(X_1 - x_0)$  is with probability one an integer value. Assume further that  $1 > \text{Prob}(X_1 = q) > 0$  (in particular, this implies that  $\frac{1}{d}(q - x_0)$  is an integer and that  $\Lambda''(\eta) > 0$ ). Then,*

$$\lim_{n \rightarrow \infty} J_n \mu_n(A) = \frac{\eta d}{1 - e^{-\eta d}}. \quad (2.10.244)$$

**Remarks:** (a). Recall that  $\Lambda^*(q) = \eta q - \Lambda(\eta)$  and  $\Lambda(\cdot)$  is  $C^\infty$  in some open neighborhood of  $\eta$  by the dominated convergence argument (for details see Lemma 2.2.1 and exercise 2.2.5).

(b). Actually the limit relations (2.10.243) and (2.10.244) remain valid even for small intervals of size of  $O(\frac{\log n}{n})$  (see exercise 2.10.1).

(c). The proof of this theorem is based on an exponential translation of a local Central Limit Theorem. This approach is applicable for the dependent case of Section 2.3 and to certain extent applies also in  $\mathbb{R}^d$ ,  $d > 1$ .

**Proof:** (a). Consider the probability measure  $\bar{\mu}$  defined by  $d\bar{\mu}(x) \triangleq e^{\eta x - \Lambda(\eta)} d\mu(x)$  and let  $Y_i \triangleq (X_i - q)/\sqrt{\Lambda''(\eta)}$ , for  $i = 1, 2, \dots, n$ . Note that  $Y_1, \dots, Y_n$  are i.i.d. random variables with  $E_{\bar{\mu}}[Y_1] = 0$  and  $E_{\bar{\mu}}[Y_1^2] = 1$  (this can be easily checked by computing the first two moments of  $X_1$  under  $\bar{\mu}$ ). Let  $F_n(x)$  denotes the distribution function of  $W_n \triangleq \frac{1}{\sqrt{n}} \sum_{i=1}^n Y_i$  under the measure  $\bar{\mu}$ .

Since  $X_i$  are non-lattice, the Berry-Esseen expansion of  $F_n(x)$  results with:

$$\lim_{n \rightarrow \infty} \left\{ \sqrt{n} \sup_x \left| F_n(x) - \Phi(x) - \frac{m_3}{6\sqrt{n}} (1 - x^2) \phi(x) \right| \right\} = 0, \quad (2.10.245)$$

where  $m_3 \triangleq E_{\bar{\mu}}[Y_1^3] < \infty$ ,  $\phi(x) = \frac{1}{\sqrt{2\pi}} \exp(-x^2/2)$  is the standard Normal density, and  $\Phi(x) = \int_{-\infty}^x \phi(\eta) d\eta$  (for the derivation of (2.10.245) see [12] page 512).

Now,

$$\begin{aligned} \mu_n(A) &= \mu_n([q, \infty)) = E_{\bar{\mu}}[e^{-n[\eta \hat{S}_n - \Lambda(\eta)]} 1_{\hat{S}_n \geq q}] = \\ &= e^{-n\Lambda^*(q)} E_{\bar{\mu}} \left[ e^{-\eta \sqrt{n\Lambda''(\eta)} W_n} 1_{W_n \geq 0} \right] = e^{-n\Lambda^*(q)} \int_0^\infty e^{-\eta \sqrt{n\Lambda''(\eta)} x} dF_n(x) \end{aligned} \quad (2.10.246)$$

since  $\hat{S}_n = q + \sqrt{\frac{\Lambda''(\eta)}{n}} W_n$ . Let  $\psi_n \triangleq \eta \sqrt{n\Lambda''(\eta)}$ . By an integration by parts in (2.10.246) one obtains

$$J_n \mu_n(A) = \sqrt{2\pi} \int_0^\infty \psi_n^2 e^{-\psi_n x} [F_n(x) - F_n(0)] dx = \sqrt{2\pi} \int_0^\infty \psi_n e^{-t} \left[ F_n\left(\frac{t}{\psi_n}\right) - F_n(0) \right] dt \quad (2.10.247)$$

Consider now

$$c_n \triangleq \sqrt{2\pi} \int_0^\infty \psi_n e^{-t} \left[ \Phi\left(\frac{t}{\psi_n}\right) + \frac{m_3}{6\sqrt{n}} \left[ 1 - \left(\frac{t}{\psi_n}\right)^2 \right] \phi\left(\frac{t}{\psi_n}\right) - \Phi(0) - \frac{m_3}{6\sqrt{n}} \phi(0) \right] dt \quad (2.10.248)$$

Comparing (2.10.247) and (2.10.248), observe that the Berry-Esseen expansion (2.10.245) yields the relation  $\lim_{n \rightarrow \infty} |J_n \mu_n(A) - c_n| = 0$ . Moreover, since

$$\sup_{x \geq 0} |\phi'(x)| < \infty, \quad \lim_{x \rightarrow 0} |\phi'(x)| = 0 \quad (2.10.249)$$

it follows by a Taylor expansion of  $\Phi\left(\frac{t}{\psi_n}\right)$  and the dominated convergence theorem that

$$\begin{aligned} \lim_{n \rightarrow \infty} c_n &= \lim_{n \rightarrow \infty} \sqrt{2\pi} \int_0^\infty \psi_n e^{-t} \left[ \Phi\left(\frac{t}{\psi_n}\right) - \Phi(0) \right] dt = \\ &= \lim_{n \rightarrow \infty} \sqrt{2\pi} \int_0^\infty e^{-t} \phi\left(\frac{t}{\psi_n}\right) dt = \sqrt{2\pi} \phi(0) \int_0^\infty e^{-t} dt = 1. \end{aligned} \quad (2.10.250)$$

This completes the proof for the non-lattice case.

(b). In the lattice case (where the range of  $Y_i$  is  $\left\{m \frac{d}{\sqrt{\Lambda''(\eta)}}\right\}_{m=-\infty}^{\infty}$ ) the Berry-Esseen expansion (2.10.245) is modified to

$$\lim_{n \rightarrow \infty} \left\{ \sqrt{n} \sup_x \left| F_n(x) - \Phi(x) - \frac{m_3}{6\sqrt{n}} (1-x^2) \phi(x) - \phi(x) g\left(x, \frac{d}{\sqrt{\Lambda''(\eta)n}}\right) \right| \right\} = 0 \quad (2.10.251)$$

where  $g(x, h) = \frac{h}{2} - (x \bmod h)$  if  $(x \bmod h) \neq 0$  and  $g(x, h) = -\frac{h}{2}$  if  $(x \bmod h) = 0$  (see [12], page 513 or [22], page 171, Theorem 6). Thus, by adopting the argument above for the lattice case, one obtains

$$\lim_{n \rightarrow \infty} J_n \mu_n(A) = 1 + \lim_{n \rightarrow \infty} \sqrt{2\pi} \int_0^\infty \psi_n e^{-t} \left\{ \phi\left(\frac{t}{\psi_n}\right) g\left(\frac{t}{\psi_n}, \frac{\eta d}{\psi_n}\right) - \phi(0) g\left(0, \frac{\eta d}{\psi_n}\right) \right\} dt \quad (2.10.252)$$

Since  $\psi_n g\left(\frac{t}{\psi_n}, \frac{\eta d}{\psi_n}\right) = g(t, \eta d)$  it follows that

$$\begin{aligned} \lim_{n \rightarrow \infty} J_n \mu_n(A) &= 1 + \lim_{n \rightarrow \infty} \sqrt{2\pi} \int_0^\infty e^{-t} \left\{ \phi\left(\frac{t}{\psi_n}\right) g(t, \eta d) - \phi(0) g(0, \eta d) \right\} dt \\ &= 1 + \sqrt{2\pi} \phi(0) \int_0^\infty e^{-t} [g(t, \eta d) - g(0, \eta d)] dt. \end{aligned} \quad (2.10.253)$$

The proof is completed by combining (2.10.253) with

$$\int_0^\infty e^{-t} [g(t, \eta d) - g(0, \eta d)] dt = \left\{ \sum_{n=0}^\infty e^{-n\eta d} \right\} \int_0^{\eta d} e^{-t} (\eta d - t) dt = \frac{\eta d}{1 - e^{-\eta d}} - 1. \quad (2.10.254)$$

□

### Exercises:

**2.10.1** (a). Let  $A = [q, q + \frac{a}{n})$ , where in the lattice case  $\frac{a}{d}$  is restricted to be an integer. Prove that for any  $a \in (0, \infty)$ , both (2.10.243) and (2.10.244) hold with  $J_n = \eta \sqrt{\Lambda''(\eta) 2\pi n} e^{n\Lambda^*(q)} \frac{1}{1 - e^{-\eta a}}$ .

(b). As a consequence of part (a) above conclude that for any set  $A = [q, q + b_n)$  both (2.10.243) and (2.10.244) hold for  $J_n$  as given in Theorem 2.10.2 as long as  $\lim_{n \rightarrow \infty} nb_n = \infty$ .

**2.10.2** (a). Let  $\eta > 0$  denote the minimizer of  $\Lambda(\lambda)$  and suppose that  $\Lambda(\lambda) < \infty$  in some open interval around  $\eta$ . Based on exercise 2.10.1, deduce that the limiting distribution of  $S_n$  conditional upon  $S_n \geq 0$  is  $\text{Exponential}(\eta)$  when  $X_1$  has a non-lattice distribution.

(b). Suppose now that  $X_1$  has a lattice distribution of span  $d$  and  $1 > \text{Prob}(X_1 = 0) > 0$ . Deduce now that the limiting distribution of  $\frac{1}{d} S_n$  conditional upon  $S_n \geq 0$  is  $\text{Geometric}(p)$  with  $p = 1 - e^{-\eta d}$  (i.e.,  $\text{Prob}(S_n = kd | S_n \geq 0) \rightarrow pq^k$  for  $k = 0, 1, 2, \dots$ ).

**2.10.3** Consider a Neyman-Pearson test with constant threshold  $\gamma \in (\bar{x}_0, \bar{x}_1)$  (see Section 2.7 for details). Suppose that  $X_1 = \log \frac{d\mu_1}{d\mu_0}(Y_1)$  has a non-lattice distribution. Let  $\lambda_\gamma \in (0, 1)$  be the unique solution of  $\Lambda'_0(\lambda) = \gamma$ . Deduce from (2.10.243) that

$$\lim_{n \rightarrow \infty} \left\{ \alpha_n e^{n\Lambda^*(\gamma)} \lambda_\gamma \sqrt{\Lambda''(\lambda_\gamma) 2\pi n} \right\} = 1 \quad (2.10.255)$$

and

$$\lim_{n \rightarrow \infty} \left\{ \frac{e^{n\gamma} \alpha_n}{\beta_n} \right\} = \frac{1 - \lambda_\gamma}{\lambda_\gamma}. \quad (2.10.256)$$

## Chapter 3

# Historical notes and references

Although much of the credit for the modern theory of large deviations and its various applications must go to Donsker and Varadhan, the topic is much older and references to the various aspects of it may be traced back to the early 1900's. Due to our own ignorance we will necessarily confine ourselves here to an incomplete list of references and historical credits. We hope to expand and correct this list at a later stage, and apologize to those who are not given due credit.

### 3.1 Chapter 2

The early development of large deviation bounds did not follow the order of our presentation. Statisticians, starting with Khinchin [18], have analysed various forms of Cramer's theorem for special random variables. See [27], [20] and [21] for additional references on this early work.

The first statement of Cramer's Theorem for distributions on  $R$  possessing densities is due to Cramer [7], who introduced the change of measure argument to this context. An extension to general distributions was done by Chernoff [6], who introduced the upper bound which was to carry his name. There exists a large body of literature concerning the applications of Cramer's theorem to the analysis of statistical tests, on which we keep silent.

Although Stirling's formula, which is at the heart of the combinatorial estimates of section 2.1, dates back at least to the 19-th century, the notion of types and bounds of the form of Lemmas 2.1.1-2.1.4 had to wait until information theorists discovered that they are useful tools for analyzing



the efficiency of codes. For early references we refer the reader to the excellent book of Gallager [12]. Our treatment of the source coding theorem in section 2.9 is a combination of the method of that book with the particular case treated by Bucklew in [4].

The credit for the extension of Cramer's theorem to the dependent case should definitely go to Gärtner [14] who considered the case in which  $\mathcal{D}_\Lambda = \mathbb{R}^d$ . Ellis [10] extended this result to the steep set up and the formulation of section 2.3 is nothing but an embellishment of his results.

The large deviations statements for Markov chains have a long history which is partially described in the historical notes of Chapter ?? . The approach taken here is based in part on the ideas of Ellis [10].

The material in section 2.5 is a large deviations proof of the results in [22].

Gibb's conditioning principle has served as a driving force behind Ruelle and Landford's treatment of large deviations (without calling it by that name), [24], [25], [19]. The form of the Gibb's principle here was proved using large deviations methods (via the method of types) by Campenhout and Cover [5] and in greater generality by Csizar [8] and Stroock-Zeitouni [28].

No references yet for Stein's lemma.

The generalized maximum likelihood of section 2.8 was considered by Hoeffding [16], whose approach we basically follow-here. The extension to general state space presented in section ?? is due to Zeitouni and Gutman [30].

Finally, the refinements of the large deviations principles discussed in sections 2.10 and ?? follow [1], [21], although some of the methods are much older and may be found in Feller's book [11].

# Bibliography

- [1] R.R. Bahadur and R. Ranga Rao. On deviations of the sample mean. *Ann. Math Statistics*, 38:1015–1027, 1960.
- [2] T. Berger. *Rate Distortion Theory*. Prentice Hall, 1971.
- [3] P. Billingsley. *Convergence of Probability Measures*. Wiley, 1968.
- [4] J. A. Bucklew. *Large deviations techniques in decision, simulation, and estimation*. Springer, 1991.
- [5] J.M. Van Campenhout and T.M. Cover. Maximum entropy and conditional probability. *IEEE Transc. Inf. Theory*, IT-27:483–489, 1981.
- [6] H. Chernoff. A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *Ann. Math., Statist.*, 23:493–507, 1952.
- [7] H. Cramèr. Sur un nouveau théorème-limite de la théorie des probabilités. In *Actualités Scientifiques et Industrielles*, volume 3 of *Colloque consacré à la théorie des probabilités*, pages 5–23. Hermann, Paris, 1938.
- [8] I. Csiszàr. I-divergence geometry of probability distributions and minimizations problems. *Ann. Probab.*, 3:146–158, 1975.
- [9] I. Ekeland and R. Temam. *Convex Analysis and Variational Problems*. North Holland, Amsterdam, 1976.
- [10] R. S. Ellis. Large deviations for a general class of random vectors. *Ann. Probab.*, 12:1–12, 1984.

- [11] W. Feller. *An Introduction to Probability Theory and Its Applications*, volume I. John Wiley and Sons, New York, 1957.
- [12] R.G. Gallager. *Information Theory and Reliable Communication*. Wiley, New York, 1968.
- [13] F.G. Gantmacher. *The Theory of Matrices*. Chelsea, 1959.
- [14] J. Gärtner. On large deviations from the invariant measure. *Theory Probab. Appl.*, 22:24–39, 1977.
- [15] U. Grenander and G. Szego. *Toeplitz Forms and their Applications*. University of California Press, 1958.
- [16] W. Hoeffding. On probabilities of large deviations. In *Proceedings of the Fifth Berkeley Symposium on mathematical Statistics and Probability*, pages 203–219. Univ. of California Press, 1965.
- [17] I. Csizar, T.M. Cover and B.S. Choi. Conditional limit theorems under Markov conditioning. *IEEE Trans. Inf. theory*, IT-33:788–801, 1987.
- [18] A.I. Khinchin. Über einen neuen grenzwertsatz der wahrscheinlichkeitsrechnung. *Math Annalen*, 101:745–752, 1929.
- [19] O.E. Landford. Entropy and equilibrium states in classical statistical mechanics. In A. Lenard, editor, *Statistical Mechanics and Mathematical problems*, volume 20 of *Lecture Notes in Physics*, pages 1–113. Springer, Berlin, 1973.
- [20] Y. V. Linnik. On the probability of large deviations for the sums of independent variables. In *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability*, pages 289–306, Berkeley, 1961. Univ. of California Press.
- [21] V. V. Petrov. *Sums of independent random variables*. Springer, Berlin, 1975. Translated by A. A. Brown.
- [22] R. Arratia, L. Gordon and M.S. Waterman. The Erdos-Renyi law in distribution for coin tossing and sequence matching. *Ann. of Statistics*, 1990.

- [23] R. T. Rockafellar. *Convex Analysis*. Princeton University Press, Princeton, 1970.
- [24] D. Ruelle. Correlation functionals. *J. math. Physics*, 6:201–220, 1965.
- [25] D. Ruelle. A variational formulation of equilibrium statistical mechanics and the Gibbs phase rule. *Comm. math. Phys.*, 5:324–329, 1967.
- [26] E. Seneta. *Non Negative Matrices and Markov Chains*. Springer-Verlag, 1981.
- [27] N. Smirnov. Über warhscheinlichkeiten grosser abweichungen. *Rec. Sco. Math. Moscou*, 40:441–455, 1933.
- [28] D.W. Stroock and O. Zeitouni. Micro canonical distributions, Gibbs' states, and the equivalence of ensembles. In R. Durrett and H. Kesten, editors, *Festschrift in honour of F. Spitzer*. Birkhauser, 1990.
- [29] T.M. Cover and J.B. Thomas. *Elements of Information Theory*. forthcoming, 1991.
- [30] O. Zeitouni and M. Gutman. On universal hypotheses testing via large deviations. *IEEE Trans. Inf. Theory*, 1991.